

ARIZONA DEPARTMENT OF TRANSPORTATION

REPORT NUMBER: FHWA-AZ99-462

RHODES - ITMS CORRIDOR CONTROL PROJECT

Final Report

Prepared by:

Douglas Gettman

Larry Head

Pitu Mirchandani

Systems and Industrial Engineering Department

University of Arizona

Tucson, Arizona 85721

May 1999

Prepared for:

Arizona Department of Transportation

206 South 17th Avenue

Phoenix, Arizona 85007

in cooperation with

U.S. Department of Transportation

Federal Highway Administration

The contents of this report reflect the views of the authors who are responsible for the facts and the accuracy of the data presented herein. The contents do not necessarily reflect the official views or policies of the Arizona Department of Transportation or the Federal Highways Administration. This report does not constitute a standard, specification, or regulation. Trade or manufacturer's names which may appear herein are cited only because they are considered essential to the objectives of the report. The U.S. Government and the State of Arizona do not endorse products or manufacturers.

Technical Report Documentation Page

1. Report No. FHWA-AZ99-462		2. Government Accession No.		3. Recipient's Catalog No.	
4. Title and Subtitle RHODES-ITMS CORRIDOR CONTROL PROJECT				5. Report Date MAY 1999	
				6. Performing Organization Code	
7. Author Douglas Gettman, Larry Head, and Pitu Mirchandani				8. Performing Organization Report No.	
9. Performing Organization Name and Address Systems and Industrial Engineering Department The University of Arizona Tucson, Arizona 85721				10. Work Unit No.	
				11. Contract or Grant No. SPR-PL-1(51)462	
12. Sponsoring Agency Name and Address ARIZONA DEPARTMENT OF TRANSPORTATION 206 S. 17TH AVENUE PHOENIX, ARIZONA 85007 ADOT Project Manager: Stephen R. Owen, P.E.				13. Type of Report & Period Covered Final Report 9/96 -5/99	
				14. Sponsoring Agency Code	
15. Supplementary Notes Prepared in cooperation with the U.S. Department of Transportation, Federal Highway Administration					
16. Abstract <p>The <i>RHODES-ITMS Corridor Control</i> project addresses real-time control of ramp meters of a freeway segment, with consideration of the traffic volumes entering and leaving the freeway from/to arterials, and the regulation of these volumes via real-time setting of ramp metering rates.</p> <p>Current approaches to traffic-responsive control of ramp meters include (a) time-of day control, (b) locally responsive strategies and (c) area-wide <i>linear programming</i> based approaches (currently implemented in parts of Europe). None of these approaches are both real-time responsive to traffic conditions <u>and</u> consider the multiple objectives of minimizing freeway travel times and decreasing congestion/queues at the interchanges.</p> <p>A control system, MILOS (Multi-objective Integrated Large-scale Optimized ramp metering System), was developed that determines ramp metering rates based on observed traffic on the freeway and its interchange arterials. MILOS temporally and spatially decomposes the ramp-metering control problem into three hierarchical subproblems: (1) monitoring and detection of traffic anomalies (to schedule optimizations at the lower levels of the control hierarchy), (2) optimization to obtain area-wide coordinated metering rates, and (3) real-time regulation of metering rates to adjust for local conditions.</p> <p>Simulation experiments were performed to evaluate the MILOS hierarchical system against (a) "no control" (i.e., when no ramp metering is in effect), (b) a locally traffic-responsive metering policy, and (c) an area-wide LP optimization problem re-solved in 5-minute intervals. Three test scenarios were simulated (1) a short "burst" of heavy-volume flows to all ramps, (2) a three-hour commuting peak, and (3) a three-hour commuting peak with a 30-minute incident occurring somewhere in the middle of the corridor. The performance results indicate that MILOS is able to reduce freeway travel time, increase freeway average speed, and improve recovery performance of the system when flow conditions become congested due to an incident.</p>					
17. Key Words Real-time traffic-adaptive control; Adaptive Ramp-metering; Freeway Traffic Control; Hierarchical Control; Simulation; Optimization.		18. Distribution Statement Document is available to the U.S. Public through the National Technical Information Service, Springfield, Virginia, 22161		23. Registrant's Seal	
19. Security Classification Unclassified	20. Security Classification Unclassified	21. No. of Pages 208	22. Price		

SI* (MODERN METRIC) CONVERSION FACTORS

APPROXIMATE CONVERSIONS TO SI UNITS

Symbol	When You Know	Multiply By	To Find	Symbol
--------	---------------	-------------	---------	--------

LENGTH

in	inches	25.4	millimeters	mm
ft	feet	0.305	meters	m
yd	yards	0.914	meters	m
mi	miles	1.61	kilometers	km

AREA

in ²	square inches	645.2	millimeters squared	mm ²
ft ²	square feet	0.093	meters squared	m ²
yd ²	square yards	0.836	meters squared	m ²
ac	acres	0.405	hectares	ha
mi ²	square miles	2.59	kilometers squared	km ²

VOLUME

fl oz	fluid ounces	29.57	milliliters	mL
gal	gallons	3.785	liters	L
ft ³	cubic feet	0.028	meters cubed	m ³
yd ³	cubic yards	0.765	meters cubed	m ³

NOTE: Volumes greater than 1000 L shall be shown in m³.

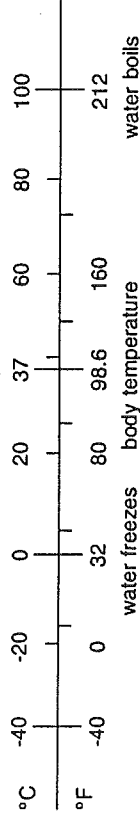
MASS

oz	ounces	28.35	grams	g
lb	pounds	0.454	kilograms	kg
T	short tons (2000 lb)	0.907	megagrams	Mg

TEMPERATURE (exact)

Symbol	When You Know	Do The Following	To Find	Symbol
--------	---------------	------------------	---------	--------

°F	Fahrenheit temperature	°F - 32 ÷ 1.8	Celcius temperature	°C
----	------------------------	---------------	---------------------	----



APPROXIMATE CONVERSIONS FROM SI UNITS

Symbol	When You Know	Multiply By	To Find	Symbol
--------	---------------	-------------	---------	--------

LENGTH

mm	millimeters	0.039	inches	in
m	meters	3.28	feet	ft
m	meters	1.09	yards	yd
km	kilometers	0.621	miles	mi

AREA

mm ²	millimeters squared	0.0016	square inches	in ²
m ²	meters squared	10.764	square feet	ft ²
m ²	meters squared	1.19	square yards	yd ²
ha	hectares	2.47	acres	ac
km ²	kilometers squared	0.386	square miles	mi ²

VOLUME

mL	milliliters	0.034	fluid ounces	fl oz
L	liters	0.264	gallons	gal
m ³	meters cubed	35.315	cubic feet	ft ³
m ³	meters cubed	1.31	cubic yards	yd ³

MASS

g	grams	0.035	ounces	oz
kg	kilograms	2.205	pounds	lb
Mg	megagrams	1.102	short tons (2000 lb)	T

TEMPERATURE (exact)

Symbol	When You Know	Do The Following	To Find	Symbol
--------	---------------	------------------	---------	--------

°C	Celcius temperature	°C x 1.8 + 32	Fahrenheit temperature	°F
----	---------------------	---------------	------------------------	----

METER: a little longer than a yard (about 1.1 yards)
LITER: a little larger than a quart (about 1.06 quarts)
GRAM: a little more than the weight of a paper clip
MILLIMETER: diameter of a paper clip wire
KILOMETER: somewhat further than 1/2 mile (about 0.6 mile)

RHODES-ITMS Corridor Control Project

PREFACE

This report documents the work performed on the *Corridor Control* Subproject of the *RHODES-ITMS* Project. This research effort was funded by the Arizona Department of Transportation (ADOT) and the Maricopa Association of Governments (MAG). Essentially, the scope of the Project was to develop a method to optimally control, in real time, the ramp meters on a segment of a freeway. The control architecture used was based on extensions of the hierarchical control concepts developed for the surface street network in the previous *RHODES* Project funded by ADOT and the Pima Association of Governments.

The *Corridor Control* subproject was directed by the principal investigators, **Pitu B. Mirchandani** and **Larry Head**, both of the Systems and Industrial Engineering Department at the University of Arizona. This report is largely based on the dissertation written by Dr. **Douglas Gettman** whose Ph.D. research was supported by the project.

In addition, Drs. Gettman, Head, and Mirchandani wish to acknowledge their appreciation to the Project's **Technical Advisory Committee (TAC)** whose continual active participation, technical input and support resulted in *MILOS* (the real-time ramp-metering system described in this report) being very relevant to freeway ramp-metering control. The following individuals served on the TAC at various times:

Jim Decker	Traffic Operations, City of Tempe
Tim Wolfe	ADOT Technology Group
Dan Powell	ADOT District 1
Tom Parlante	ADOT Traffic Engineering
Glenn Jonas	ADOT Freeway Management
Jim Shea	ADOT Traffic Engineering
Sarath Joshua	Maricopa Association of Governments (previously at ATRC, ADOT)
Paul Ward	Maricopa Association of Governments
Roy Turner	Maricopa Association of Governments
Pierre Pretorius	Maricopa County Transportation and Development Agency
Don Wiltshire	Maricopa County Transportation and Development Agency
Alan Hansen	Federal Highway Administration
Tom Fowler	Federal Highway Administration
Jessie Yung	Federal Highway Administration
Steve Owen	RHODE-ITMS Project Manager, Arizona Transportation Research Center (ATRC), ADOT

The contents of this report reflect the views of the authors who are responsible for the facts and the accuracy of the data presented herein. The contents do not necessarily reflect the official views of the Arizona Department of Transportation, Maricopa Association of Governments or the Federal Highway Administration. This report does not constitute a standard, specification or regulation.

RHODES-ITMS Corridor Control Project

EXECUTIVE SUMMARY

The RHODES-Integrated Traffic Management System (ITMS) Program addresses the design and development of a real-time traffic adaptive control system for an integrated system of freeways and arterials. The overall program was initiated in December 1993, jointly funded by the Arizona Department of Transportation (ADOT) through the State Planning and Research Program budget and the Maricopa Association of Governments (MAG). The *RHODES-ITMS* program is overseen by ADOT's Arizona Transportation Research Center.

Subsequently, in September 1996, the *RHODES-ITMS Corridor Control* research project was initiated which specifically addresses **real-time control of ramp meters** of a freeway segment, with consideration of the traffic volumes entering and leaving the freeway from/to arterials, and the regulation of these volumes via real-time setting of ramp metering rates. This is the final report for the *RHODES-ITMS Corridor Control Project*.

RESEARCH CONCEPTS

Current approaches to controlling ramp meters to respond to varying traffic conditions, those reviewed in this report, include (a) time-of day control, (b) locally responsive strategies (one such strategy is currently under consideration by ADOT), and (c) area-wide *linear programming* (LP) based approaches (currently implemented in parts of Europe). None of these approaches are both fully responsive in real-time to prevailing and predicted traffic conditions and consider the multiple objectives of minimizing freeway travel times and decreasing congestion/queues at the interchange ramps and the corresponding arterial intersections.

In this research, a control system was developed, referred to as MILOS (Multi-objective Integrated Large-scale Optimized ramp metering System), that determines ramp metering rates based on observed and predicted traffic on the freeway and its interchange arterials. MILOS has an hierarchical architecture to address the complexities of the real-time freeway management problem, namely, (a) the dynamic and stochastic nature of state changes, and (b) the existence of multiple objectives. MILOS temporally and spatially decomposes the ramp-metering control problem into three hierarchical subproblems: (1) monitoring and detection of traffic anomalies (to schedule optimizations at the lower levels of the control hierarchy), (2) optimization to obtain area-wide coordinated metering rates, and (3) real-time regulation of metering rates to adjust for local conditions.

The area-wide coordination problem at the second level of the hierarchical control system is modeled as a "*quadratic programming*" (QP) optimization problem that considers the impact of queue growth on the adjacent interchanges. A multi-criterion objective function is used to trade-off between freeway travel times and congestion/queues at the interchanges. The resulting nominal solution of the second-level area-wide optimization

problem is then provided to the third-level control function which locally adjusts these nominal ramp-meter rates.

The third-level problem, referred to as predictive-cooperative real-time (PC-RT) rate regulation problem, modifies, if necessary, the ramp metering rates based on local traffic at each interchange. The PC-RT algorithm is based on a linear programming formulation that uses a linearized approximation of a macroscopic freeway flow model (in terms of dynamic difference equations). The PC-RT algorithm pro-actively utilizes opportunities to disperse queues or hold back additional vehicles when freeway and ramp traffic conditions are appropriate. The cost coefficients of the LP optimization objectives are based on the multi-objectives trade-offs considered in the second-level area-wide coordination problem.

The optimization runs of the area-wide coordination problem and the PC-RT rate regulation problem at each ramp are scheduled for execution by the highest-level ramp-demand/freeway-flow monitoring system that is based on concepts from “*statistical process control*” in production systems. Basically, this system functions as follows: When the monitored conditions are within the expected variances in ramp demands and freeway flows, no optimization run is scheduled to obtain new ramp metering rates; when the conditions are outside the expected variances then either the PC-RT algorithm (LP) is run if the deviations are not too large, or the area wide QP is run when the deviations are large, to obtain new ramp metering rates.

RESULTS

Simulation experiments were performed to evaluate the MILOS hierarchical system against (a) “no control” (i.e., when no ramp metering is in effect), (b) a locally traffic-responsive metering policy currently under consideration by ADOT, and (c) an area-wide LP optimization problem re-solved in 5-minute intervals. The simulation model was of a small freeway corridor in metropolitan Phoenix, Arizona (seven miles of State Route 202 with 7 off-ramps, 4 controllable on-ramps, and one freeway-freeway on-ramp (hypothesized as controllable). Three test scenarios were simulated (1) a short “burst” of heavy-volume flows to all ramps, (2) a three-hour commuting peak, and (3) a three-hour commuting peak with a 30-minute incident occurring somewhere in the middle of the corridor.

The performance results indicate that MILOS is able to reduce freeway travel time, increase freeway average speed, and improve recovery performance of the system when flow conditions become congested due to an incident. Specifically, when comparing with the “no control” case, freeway travel times were lowered by 8% - 36%, speeds were increased by 3% - 18%, and recovery times were reduced by 6% - 25%. It also performed better than the area-wide LP optimization that has been reported to perform well in Europe. Locally responsive strategy performed well in light to moderate traffic volumes, and, in fact, had lower freeway travel times and faster speeds than MILOS; however, for heavier volumes and incidents it had larger ramp-queues and longer recovery times than MILOS.

AREAS OF FUTURE WORK

This research project identified several interesting future research, development and deployment efforts. Development of (1) an algorithm to decompose a region into subnetworks for MILOS control, (2) a model on route diversion and (3) methods to estimate ramp demands and interchange turning probabilities are promising research areas. Integration of MILOS with traffic-adaptive surface-street signal control and incident/anomaly detection systems are developmental efforts that could make traffic management even more real-time responsive. Finally, field testing and the deployment of MILOS (and its future enhancements) should be an on-going effort towards the ITS goal of implementing advanced traffic management systems that are safer, more efficient and beneficial to the traveling public.

ACKNOWLEDGMENTS

The *RHODES-ITMS Corridor Control Project* was completed in January 1999. Project oversight was provided by a Technical Advisory Committee (TAC) comprising of representatives from key agencies. The project was administered by the Arizona Transportation Research Center of ADOT. The following individuals served on the TAC at various times:

Steve Owen	RHODES-ITMS Project Manager, Arizona Transportation Research Center (ATRC), ADOT
Tim Wolfe	ADOT Technology Group
Dan Powell	ADOT District 1
Tom Parlante	ADOT Traffic Engineering
Glenn Jonas	ADOT Freeway Management
Jim Shea	ADOT Traffic Engineering
Alan Hansen	Federal Highway Administration
Tom Fowler	Federal Highway Administration
Jessie Yung	Federal Highway Administration
Sarath Joshua	MAG (previously at ATRC, ADOT)
Paul Ward	MAG
Roy Turner	MAG
Pierre Pretorius	Maricopa County Transportation and Development Agency
Don Wiltshire	Maricopa County Transportation and Development Agency
Jim Decker	Traffic Operations, City of Tempe

Table of Contents

CHAPTER 1: PROBLEM OVERVIEW	1
INTRODUCTION.....	1
THE FUNDAMENTAL FREEWAY MANAGEMENT PROBLEM.....	2
ISSUES IN APPLICATION OF RAMP METERING AS A METHOD OF FREEWAY MANAGEMENT	4
THE MULTI-OBJECTIVE APPROACH TO FREEWAY SYSTEM MANAGEMENT	6
RESEARCH METHODOLOGY	7
SUMMARY OF THE FORTHCOMING CHAPTERS	8
CHAPTER 2: RAMP METERING LITERATURE REVIEW.....	10
COSTS AND BENEFITS OF RAMP METERING	10
TYPES OF RAMP METERING ALGORITHMS	11
TIME-OF-DAY METERING ALGORITHMS	11
LOCAL TRAFFIC-RESPONSIVE RAMP METERING ALGORITHMS	11
HYBRID RAMP METERING CONTROL ALGORITHMS	12
INTEGRATED FREEWAY/SURFACE-STREET METERING ALGORITHMS	14
SUMMARY.....	14
CHAPTER 3: HIERARCHICAL RAMP METERING CONTROL SYSTEM STRUCTURE	16
INTRODUCTION.....	16
MULTI-LEVEL METHODS IN HIERARCHICAL CONTROL.....	16
MULTI-LAYER HIERARCHICAL CONTROL SYSTEMS.....	17
SET-POINT REGULATION METHODS	18
THE MILOS HIERARCHICAL STRUCTURE	18
<i>Modal decomposition of the MILOS hierarchy.....</i>	<i>21</i>
<i>Subnetwork identification</i>	<i>22</i>
<i>SPC-based anomaly detection and optimization scheduling layer.....</i>	<i>24</i>
<i>Area-wide coordination layer</i>	<i>24</i>
<i>Predictive-cooperative real-time rate regulation layer.....</i>	<i>25</i>
INTEGRATION OF MILOS WITH NECESSARY EXTERNAL SYSTEMS	25
SUMMARY.....	26
CHAPTER 4: FREEWAY MACROSIMULATOR	27
MODEL CONSTRUCTION	27
MODELING FLOW IN HEAVY CONGESTION	31
DYNAMIC MODELING OF RAMP QUEUES	32
SUMMARY OF MACROSCOPIC MODEL	33
STOCHASTIC EFFECTS AND DIVERSION BEHAVIOR.....	35
SUMMARY.....	35
CHAPTER 5: AREA-WIDE COORDINATION PROBLEM	37
INTRODUCTION.....	37
MATHEMATICAL DESCRIPTION OF THE AREA-WIDE COORDINATION PROBLEM.....	37
<i>Derivation of the objective function</i>	<i>39</i>
<i>Consideration of queue storage limits.....</i>	<i>39</i>
<i>Development of a multi-criteria objective function.....</i>	<i>41</i>
<i>Setting costs according to interchange congestion level</i>	<i>42</i>
<i>Integration of surface-street flows in ramp demands</i>	<i>43</i>
<i>Quadratic objective function summary.....</i>	<i>44</i>
RESOLVING INFEASIBILITY	45
AREA-WIDE COORDINATION PROBLEM SUMMARY	46
OPERATION UNDER SEVERE CONGESTION.....	47
INTEGRATION WITH PREDICTIVE-COOPERATIVE REAL-TIME RATE REGULATION LAYER.....	48
PRELIMINARY EVALUATION OF AREA-WIDE COORDINATION PROBLEM ON SMALL EXAMPLE	48

<i>Influence of β</i>	50
<i>Comparison of area-wide metering rate settings in macrosimulation</i>	51
<i>Simulation test with extended queue dissipation</i>	54
SUMMARY.....	57
CHAPTER 6: PREDICTIVE-COOPERATIVE REAL-TIME RATE REGULATION ALGORITHM	58
INTRODUCTION.....	58
ADDING QUEUE MANAGEMENT TO STATE FEEDBACK CONTROL METHODS	59
CENTRAL CONCEPT OF PC-RT RATE REGULATION ALGORITHM.....	59
<i>Anticipated effects of PC-RT rate regulation algorithm</i>	61
<i>Basic function of the PC-RT rate regulation algorithm</i>	62
<i>Reasonable and important assumptions</i>	63
<i>Linearization about an equilibrium state</i>	64
<i>Elimination of the dynamic speed equation</i>	66
<i>Structure of the PC-RT objective function</i>	69
<i>Queue growth modeling</i>	70
<i>Control variable modeling</i>	72
<i>Derivation of the PC-RT cost coefficients from the QP solution</i>	72
<i>Computational procedure to obtain cost coefficients</i>	73
<i>Modification to the linearization procedure for unstable conditions</i>	74
SUMMARY OF PC-RT MATHEMATICAL FORMULATION	81
DIFFICULTY IN SOLVING THE MONOLITHIC PC-RT OPTIMIZATION PROBLEM.....	82
<i>Decomposition of full optimization problem into subproblems</i>	83
SCENARIO PREDICTION	86
<i>Scenario prediction example</i>	88
<i>Construction of scenario rate tables</i>	90
<i>Infeasible PC-RT scenarios</i>	90
SUMMARY	92
CHAPTER 7: SPC-BASED ANOMALY DETECTION	94
INTRODUCTION.....	94
OVERVIEW OF STATISTICAL PROCESS CONTROL	94
RELATIONSHIP OF SPC CONCEPTS TO FREEWAY CONTROL.....	97
JUSTIFICATION OF APPROXIMATELY-CONSTANT DEMAND.....	98
SPC COMPUTATIONAL PROCEDURE	100
TRANSITION TO A NEW \bar{X} LEVEL.....	102
OTHER ISSUES IN SPC-BASED ANOMALY DETECTION.....	104
SUMMARY.....	105
CHAPTER 8: MILOS SOFTWARE IMPLEMENTATION	106
INTRODUCTION.....	106
INITIALIZATION MODULE	108
SPC-BASED ANOMALY DETECTION MODULE	108
APPLYING NEW RATES	110
AREA-WIDE COORDINATION MODULE	110
PC-RT OPTIMIZATION MODULE.....	111
EXAMPLE OF MILOS OPERATION.....	113
SUMMARY.....	115
CHAPTER 9: SIMULATION EXPERIMENTS	117
INTRODUCTION.....	117
STRUCTURE OF THE SIMULATION EXPERIMENT.....	117
CALIBRATION OF MACROSCOPIC MODEL TO SR202 CORSIM OUTPUT.....	120
TEST CASE #1	127
<i>Results for test case #1</i>	128

TEST CASE #2	142
<i>Results for test case #2</i>	142
TEST CASE #3	157
<i>Results for test case #3</i>	158
SUMMARY	172
CHAPTER 10: CONCLUSIONS	174
GENERAL RESULTS	174
MILOS VERSUS NO CONTROL	174
MILOS VERSUS THE LOCALLY TRAFFIC-RESPONSIVE METHOD	175
MILOS VERSUS LP METHOD	176
SUMMARY	177
DIRECTIONS FOR FURTHER RESEARCH	181
APPENDIX A: SUPPORTING DATA	183
REFERENCES.....	187

List of Figures

Figure 1- 1. Empirical speed-volume measurements.....	2
Figure 1- 2. 15-minute flow time-series indicating congestion.....	3
Figure 1- 3. 15-minute speed time-series indicating congestion.....	3
Figure 2- 1. Typical hierarchical control system structure.....	13
Figure 3- 1. The pyramid structure of the MILOS hierarchy.....	20
Figure 3- 2. Example freeway network.....	22
Figure 3- 3. Initial decomposition of freeway network.....	23
Figure 3- 4. Alternative decomposition of freeway network.....	23
Figure 4- 1. Typical shape of soft-limiter function.....	33
Figure 5- 1. Ramp meter demand sources.....	44
Figure 5- 2. Example problem.....	48
Figure 5- 3. Density evolution comparison.....	53
Figure 5- 4. Queue growth comparison.....	54
Figure 5- 5. Density evolution comparisons for evaluation example 2.....	55
Figure 5- 6. Queue growth comparisons, evaluation example 2.....	56
Figure 6- 1. Prescribed maximum queue growth rate.....	71
Figure 6- 2. Example of incorrect wave-speed model for congested section.....	75
Figure 6- 3. Overcapacity segment results in upstream area-wide flow limitations.....	77
Figure 6- 4. Alternative model for the over-capacity situation.....	78
Figure 6- 5. Re-linearization for PC-RT and periodic solution of the QP.....	80
Figure 6- 6. Typical overlapping subsystem decomposition.....	84
Figure 6- 7. Predicted trends for a given subproblem.....	88
Figure 7- 1. Typical SPC control chart.....	95
Figure 7- 2. SPC limits and part specifications.....	96
Figure 7- 3. SPC chart showing sampled time-series.....	96
Figure 7- 4. Detector time-series and underlying detection history.....	97
Figure 7- 5. Re-evaluation of “approximately constant” demand level.....	98
Figure 7- 6. Jumps between approximately-constant demand levels.....	99
Figure 8- 1. MILOS operational flow chart.....	107
Figure 8- 2. SPC anomaly detection flow chart.....	109
Figure 8- 3. Area-wide coordination flow chart.....	111
Figure 8- 4. PC-RT optimization flow chart.....	113
Figure 8- 5. MILOS operational example.....	114
Figure 8- 6. MILOS operational example, continued.....	115
Figure 9- 1. State Route 202 CORSIM link-node diagram.....	121
Figure 9- 2. Comparison of density and speed measurements.....	125
Figure 9- 3. SR202 comparisons, with stochastic input flows.....	126
Figure 9- 4. Comparison of freeway travel time distributions.....	130
Figure 9- 5. Comparison of queue time distributions.....	130
Figure 9- 6. Comparison of average speed distributions.....	131
Figure 9- 7. Comparison of recovery time distributions.....	131
Figure 9- 8. Comparison of densities: no control.....	132
Figure 9- 9. Comparison of densities: TR w/QM.....	132
Figure 9- 10. Comparison of densities: LP, resolved each 5-minutes.....	133

Figure 9- 11. Comparison of densities: MILOS	133
Figure 9- 12. Comparison of queue growth: TR w/QM.....	134
Figure 9- 13. Comparison of queue growth: LP, resolved each 5-minutes.....	134
Figure 9- 14. Comparison of queue growth: MILOS.....	135
Figure 9- 15. Comparison of metering rates: no control.....	136
Figure 9- 16. Comparison of metering rates: TR w/QM.....	136
Figure 9- 17. Comparison of metering rates: LP, resolved each 5-minutes.....	137
Figure 9- 18. Comparison of metering rates: MILOS.....	137
Figure 9- 19. Total vehicles in system: no control.....	138
Figure 9- 20. Total vehicles in system: TR w/QM	139
Figure 9- 21. Total vehicles in system: LP, resolved each 5-minutes.....	140
Figure 9- 22. Total vehicles in system: MILOS	141
Figure 9- 23. Comparison of total travel time distributions	144
Figure 9- 24. Comparison of queue time distributions.....	144
Figure 9- 25. Comparison of average speed distributions.....	145
Figure 9- 26. Comparison of recovery time distributions	145
Figure 9- 27. Comparison of densities: no control.....	146
Figure 9- 28. Comparison of densities: TR w/QM	146
Figure 9- 29. Comparison of densities: LP, resolved each 5-minutes	147
Figure 9- 30. Comparison of densities: MILOS	147
Figure 9- 31. Comparison of queue growth: no control.....	148
Figure 9- 32. Comparison of queue growth: TR w/QM.....	148
Figure 9- 33. Comparison of queue growth: LP, resolved each 5-minutes.....	149
Figure 9- 34. Comparison of queue growth: MILOS.....	149
Figure 9- 35. Comparison of metering rates: no control.....	150
Figure 9- 36. Comparison of metering rates: TR w/QM.....	150
Figure 9- 37. Comparison of metering rates: LP, resolved each 5-minutes.....	151
Figure 9- 38. Comparison of metering rates: MILOS.....	151
Figure 9- 39. Total vehicles in system: no control.....	152
Figure 9- 40. Total vehicles in system: TR w/QM	153
Figure 9- 41. Total vehicles in system: LP, resolved each 5-minutes.....	154
Figure 9- 42. Total vehicles in system: MILOS	155
Figure 9- 43. SR202 model indicating incident location.....	157
Figure 9- 44. Comparison of total travel time distributions	160
Figure 9- 45. Comparison of queue time distributions.....	160
Figure 9- 46. Comparison of average speed distributions.....	161
Figure 9- 47. Comparison of recovery time distributions	161
Figure 9- 48. Comparison of densities: no control.....	162
Figure 9- 49. Comparison of densities: TR w/QM	162
Figure 9- 50. Comparison of densities: LP, resolved each 5-minutes	163
Figure 9- 51. Comparison of densities: MILOS	163
Figure 9- 52. Comparison of queue growth: no control.....	164
Figure 9- 53. Comparison of queue growth: TR w/QM.....	164
Figure 9- 54. Comparison of queue growth: LP, resolved each 5-minutes.....	165
Figure 9- 55. Comparison of queue growth: MILOS.....	165
Figure 9- 56. Comparison of metering rates: no control.....	166

Figure 9- 57. Comparison of metering rates: TR w/QM.....	166
Figure 9- 58. Comparison of metering rates: LP, resolved each 5-minutes	167
Figure 9- 59. Comparison of metering rates: MILOS.....	167
Figure 9- 60. Total vehicles in system: no control.....	168
Figure 9- 61. Total vehicles in system: TR w/QM	169
Figure 9- 62. Total vehicles in system: LP, resolved each 5-minutes.....	170
Figure 9- 63. Total vehicles in system: MILOS	171
Figure 10- 1. Comparison of typical metering rates, MILOS and LP	179
Figure 10- 2. Side-by-side comparison of metering rates for LP and MILOS	180

List of Tables

Table 1- 1. Characteristics of the proposed freeway control system	7
Table 4- 1. Parameters in macroscopic simulation equations.....	29
Table 5- 1. Route-proportional matrix of example problem.....	49
Table 5- 2. Ramp interchange data	50
Table 5- 3. Comparison of metering rate coordination methods	50
Table 5- 4. Rate comparison for $\beta=1$	51
Table 5- 5. Parameters for example problem.....	51
Table 5- 6. Initial conditions for simulation run.....	52
Table 5- 7. Preliminary method comparisons.....	54
Table 5- 8. Demand volumes for evaluation example 2	55
Table 5- 9. Initial conditions for evaluation example 2	55
Table 5- 10. Performance comparisons, evaluation example 2	56
Table 6- 1. Ramp meter demand predictions.....	89
Table 6- 2. Upstream freeway demand predictions	89
Table 6- 3. Rate table for next minute.....	90
Table 6- 4. Rate table with infeasible optimization problems	92
Table 9- 1. Traffic-responsive metering rates and thresholds	118
Table 9- 2. Route proportional matrices.....	123
Table 9- 3. Input volumes.....	123
Table 9- 4. Initial conditions and parameter values for State Route 202.....	124
Table 9- 5. State Route 202 macroscopic simulation parameters.....	126
Table 9- 6. Mean input rates, test case #1	128
Table 9- 7. Performance results of test case #1	129
Table 9- 8. Average volume rates in each time segment, test case #2.....	142
Table 9- 9. Performance comparisons, test case #2	143
Table 9- 10. Average volume rates in each time segment, test case #3	157
Table 9- 11. Performance comparisons, test case #3	159
Table 9- 12. Comparison of MILOS results with alternatives.....	173
Table 9- 13. Qualitative comparisons of MILOS versus other algorithms	173

Chapter 1: Problem overview

Introduction

Traffic delay due to congestion on freeways and surface streets was approximated at 1.2 billion vehicle-hours in the United States in 1984 and projected to reach 6.9 billion vehicle-hours by 2005 [Lindley, 1987]. User costs associated with traffic delay were estimated at \$100 billion in 1990 for the U.S. [Euler, 1990]. Total trips and commuter miles are expected to grow significantly in most metropolitan areas as the trend towards suburban sprawl continues. Construction of additional freeway lanes and wider surface streets is certainly needed to respond to such societal needs. However, in many situations, it is not possible to address capacity needs in dense urban areas with new construction. In these situations, capacity increases are possible only by adding modes of travel (rail, subway, etc.) or reducing traveler delays with more efficient management of the available system capacity.

Introduction of traffic management devices and systems (i.e. traffic signal systems, ramp metering systems, lane channelization, HOV, etc.) has been shown to reduce delays and increase capacity [FHWA, 1985]. Early estimates of the impacts of ITS technologies range from 10% reduction in emissions (relative to the projected increase in total vehicle-miles traveled) to 20% savings in vehicle-delay and 30% reduction in stops [Mobility 2000, 1989]. These results are realized *without* significant spending on road widening and adding miles of freeway. Benefit/cost ratios of 16:1 and 22:1 have been reported for investment in ITS technologies in Los Angeles and Texas, respectively [Mobility 2000, 1989].

One of the most significant contributors to total vehicle-hours of delay is the daily commute of travelers on the freeway from their homes to their place of business and vice-versa. On average, over 38% of total vehicle-hours of delay are recurrent, i.e. occurring during the commuting hours, often referred to as the *peak* periods [Lindley, 1987]. Accidents and anomalous events, sometimes referred to as *nonrecurrent* congestion-related delay, account for the majority of the remaining delay factors. Thus, the largest

impacts on the reduction of total system delay for *freeway operations*, are realized by efficiently managing the critical peak times and managing incident conditions effectively.

The fundamental freeway management problem

The central problem in freeway management can be described best by presenting the *fundamental* diagram of freeway flow. Figure 1-1 illustrates speed versus volume on a typical freeway segment in Phoenix, AZ [Technical Advisory Committee, 1997]. In this figure, the upper concentration of points represents the uncongested flow "regime", where traffic flows smoothly, i.e. without significant travel delays. The lower, less populated collection of points represents 15-minute intervals when this freeway section was congested, incurring traffic delays to travelers in and upstream of this section.

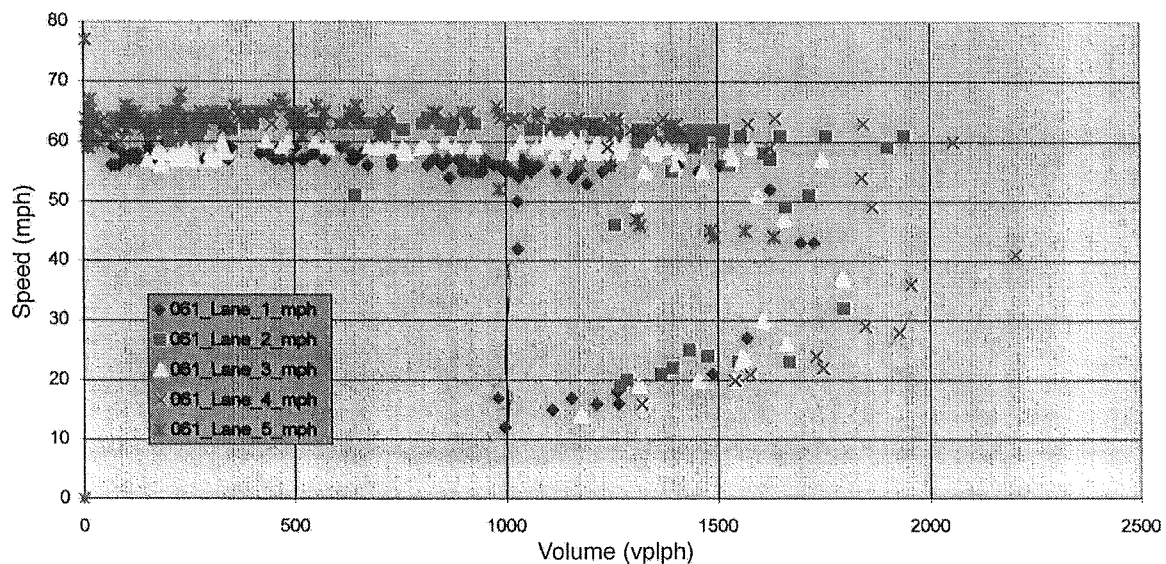


Figure 1- 1. Empirical speed-volume measurements

Figures 1-2 and 1-3 depict the time-series of the same points in Figure 1-1. These figures demonstrate that as the volume rises it becomes increasingly more precarious that the speed will drop sharply and the system will "transition" to the *congested* flow regime. As the system approaches closer and closer to the maximum flow rate, the transition can be initiated given increasingly smaller shocks [Newell, 1993]. That is to say that the higher the volume becomes, the more sensitive the system is to small anomalies in flow such as anomalies induced by merging *platoons* of vehicles.

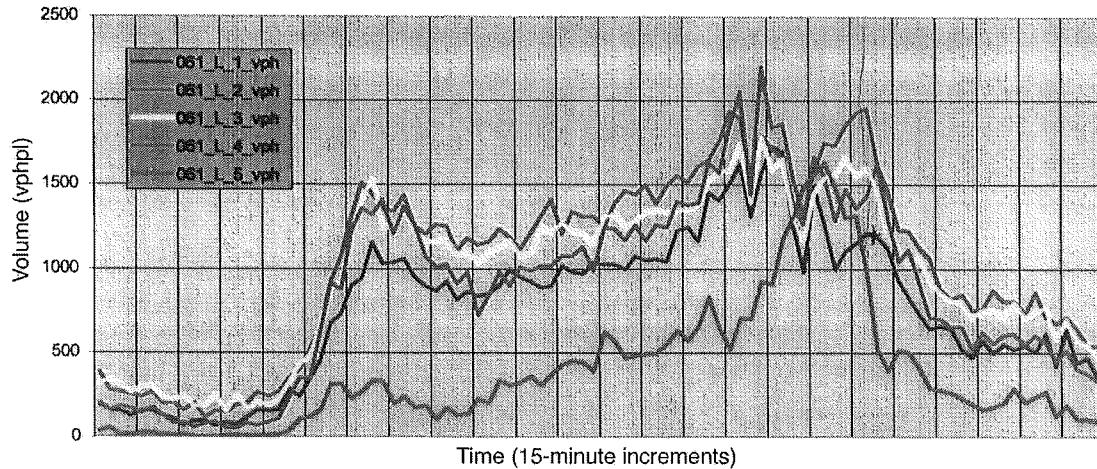


Figure 1- 2. 15-minute flow time-series indicating congestion

This increased susceptibility to shocks is due to the fact that the time between adjacent vehicles in a single lane (i.e. *headways*) *must* become smaller when the volume increases. Thus, when any interruption of the flow process occurs, drivers have less time to react to changes in the speed of the vehicle they are following and tend to over-react; braking sharply. This sharp braking can cause immediate transition to congested flow, since the volume is not reduced at the same rate as the speed.

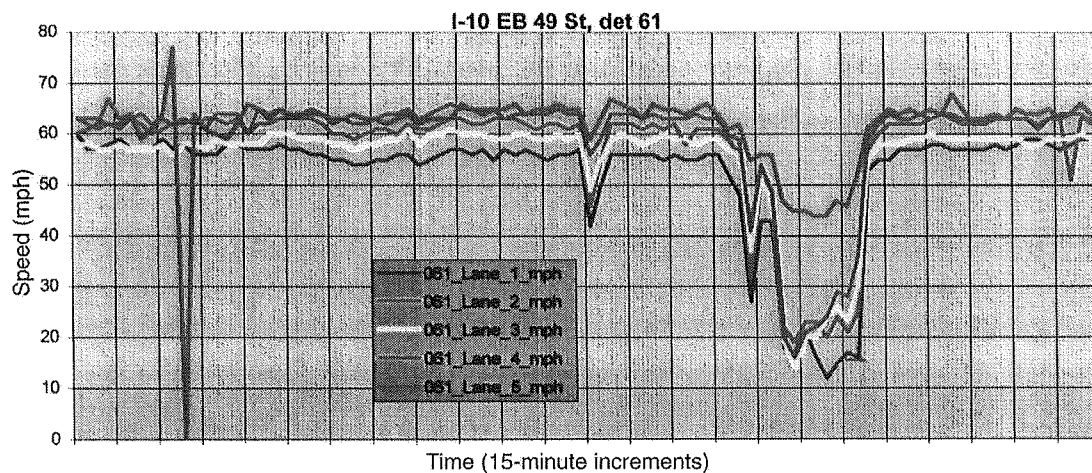


Figure 1- 3. 15-minute speed time-series indicating congestion

Given this fact, there is an obvious difficulty in operating the freeway system at the highest volumes (and best use of available capacity) near the upper end of the characteristic curve of Figure 1-1. This point is the *most susceptible* to transitioning into

the undesirable congested flow regime. The difficulty is amplified because, as indicated by Figures 1-2 and 1-3, capacity reductions usually occur when demand is greatest during the peak morning and afternoon commuting times. Thus, without some form of demand/capacity management, the sheer volume of demand for the freeway system drives the network into congestion.

There are several basic *technological* forms of freeway management available to address the freeway flow trade-off:

- (a) Advising travelers to avoid certain sections and/or change their departure times,
- (b) Advising travelers (unilaterally) to maintain a certain speed, or
- (c) Restricting access to the freeway system at certain locations.

Solution (a) describes passive advanced traveler information systems (ATIS) methods which is outside the scope of this dissertation. Solution (b) has been investigated [Karaaslan, et al., 1990; Smulders, 1993], but is likely to have compliance problems in the U.S. Automated highway systems (AHS) eliminate this compliance problem, but are relatively far from mass implementation due to regulatory concerns and cost issues [Bender, 1991]. Solution (c) is generally referred to as *ramp metering*, and is the primary topic of this research.

Issues in application of ramp metering as a method of freeway management

Ramp metering systems have existed since the early 1960's and have been used effectively in many municipalities [Carlson, 1979; Marsden, 1981; Jacobson, 1989; Haj-Salem et al., 1990; Hallenbeck and Nisbet, 1993; Wright, 1993] and in others with less conclusive benefits [Lipp et al., 1992], but there is general consensus that ramp metering systems can provide substantial benefits in throughput, travel-time, and congestion reduction when applied *appropriately*. In fact, the effectiveness of on-ramp metering has been substantial enough that a recent study proposed *main-line* metering as a tool for congestion management [Haboian, 1997].

There does *not* exist a consensus, however, of what constitutes the most *effective* metering rate(s), arising from the fact that the freeway management problem is difficult to solve to optimality. Consider the following complicating factors:

- (a) The state variables (i.e. volume, density, speed, ramp queues, etc.) change *dynamically over time*;
- (b) Although the behavior of *individual* travellers is somewhat deterministic in that the drivers know where they are going, or at least have a trip purpose, the behavior of traffic as a *stream* is stochastic and difficult to predict over long time horizons;
- (c) Flow *anomalies* (e.g. accidents, friction effects) occur at random, unpredictable intervals;
- (d) The ramp metering problem is a *multi-dimensional* one, with a large number of state and control variables;
- (e) The state variables are only *partially observable* at a limited number of fixed locations where detectors are installed, and
- (f) There are *multiple stakeholders* and, consequently, *multiple objectives* need to be addressed in any freeway management policy.

The concerns of multiple stakeholders arise when you consider the fact that the freeway system exists *embedded inside* and *interacting with* a larger network of surface-streets and other modes of transportation. Previous research has not considered the complexity associated with considering multiple stakeholders by:

- (1) assuming that the freeway system exists in virtual *isolation* from the larger surface-street network (e.g. usually assuming that ramp queuing capacity is infinite), and
- (2) choosing a single *system-optimal* optimization criterion that considers only the effects of metering decisions on freeway conditions.

Freeway management policies developed or proposed to date have included, however, considerations of dynamic state changes, stochasticity, multi-dimensionality, unpredictability, and partial-observability in the freeway management problem.

The question remains, however, whether a *system-optimal* policy for the freeway control problem *alone* is "system-optimal" for the *entire* transportation network. Here we indicate the transportation network as the entire system of freeways and surface streets in a metropolitan area or municipality (or collection thereof). In fact, it is entirely plausible that the freeway management policy "optimal" to the freeway conditions could be counterproductive to the entire transportation network because the interactions between the two surface street system and the freeway system are neglected.

This is especially true for ramp metering methods that (inevitably) create *queues* at the freeway access points. These queues, if not suitably managed, can interfere with operation of the surface-street system by extending into the adjacent interchange (commonly known as *spillback*). Thus, the objectives at the ramp interface of the surface-street manager and the freeway system manager conflict. The surface-street manager would like to keep the ramp queue as short as possible and the freeway manager would like to keep the queue as long as possible typically during congested conditions.

The multi-objective approach to freeway system management

The central issue addressed by this research is the consideration of the important interaction between the surface-street system and the freeway system and their traffic objectives in the development of a freeway access control (ramp metering) system. This problem is addressed by using a *multi-objective* solution methodology. The trade-off solutions produced by this solution method are defined by *combining* the two conflicting objectives into a single, multi-criterion objective as opposed to other methods that enumerate Pareto solutions [Haimes, et al., 1990]. In addition, because of the relative size of the carrying capacity of the freeway with respect to the adjacent surface-street system, the trade-off solution point is selected to *maintain* freeway performance that is at least as good as management policies that do *not* consider the interactions of the two sub-systems. It will be shown in Chapter 9, via simulation, that acceptable freeway performance similar to area-wide control methods that do *not* consider the effects on the interchanges can be obtained by implementing a compromise solution while, at the same time, providing queue management.

Research methodology

To mitigate the complicating factors of the multi-objective ramp queue management issue, this research uses a structured approach based on previous work in freeway ramp metering control systems, but utilizing new technologies where appropriate. Table 1-1 indicates the characteristic of the research methodology that addresses each of the complicating factors.

Complicating factor	Mitigating control system characteristic
Dynamic state changes	Rolling-horizon optimization Temporal-spatial decomposition
Stochasticity	SPC-based anomaly detection Temporal-spatial decomposition
Multi-dimensionality	Temporal-spatial decomposition
Unpredictability	Predictive scenario optimization SPC-based anomaly detection
Partial-observability	Predictive scenario optimization Rolling-horizon optimization
Multiple objectives	Multi-objective criterion functions Cost coefficient trade-off weights

Table 1- 1. Characteristics of the proposed freeway control system

In brief, Table 1-1 identifies the characteristics of a hierarchical control system that decomposes the large-scale freeway ramp metering into a series of optimization problems of varying temporal and spatial resolution. The optimization problems are re-solved as the parameters and conditions of the system change to continually adjust the control strategy to the real-time behavior of the system. In addition, to mitigate the unpredictability of the future system state, a predictive scenario-based optimization scheme is implemented in real-time to prepare the local subsystem for the next short-term stochastic disturbance.

Summary of the forthcoming chapters

The remainder of this document is structured as follows; Chapter 2 presents a brief overview of previous work on the ramp metering problem. Chapter 3 outlines the hierarchical structure of the research methodology and the temporal/spatial decomposition of the control problem. Although hierarchical treatment of freeway management is not new, the specific hierarchy proposed in this project is novel, in particular the identification of subnetworks from a large-scale freeway system, and the basis for interaction between the area-wide layer and the locally traffic-responsive layer are new. Chapter 4 presents a popular and useful model of freeway traffic flow modified slightly to more accurately represent the ramp-freeway interface under the presence of congestion. Chapter 5 presents the area-wide coordination component of the hierarchical control system that considers the impact of queue growth on the adjacent interchanges in the optimization model. This optimization model is based on models available in the literature but incorporates several additions: (1) a new multi-criterion objective function and trade-off structure, (2) an alternative treatment of queue growth constraints, and (3) modeling of demands from surface-street interchange flows.

Chapter 6 presents the locally traffic-reactive, predictive-cooperative real-time rate regulation algorithm that provides additional capacity at the freeway/surface-street interface. The basis for this optimization model is not new (i.e. linearization of the nonlinear macroscopic flow model of Chapter 4), but the formulation of the scenario-based linear-programming problem is new. The link to the solution of the area-wide coordination problem of Chapter 5 using the dual information is entirely novel.

Chapter 7 presents the statistical process control concepts used to monitor system operation and, in real-time, identify perturbations to the system states. This structure of demand estimation and fluctuation identification in the context of freeway management systems is an entirely new treatment of this modeling/estimation/optimization procedure. Chapter 8 summarizes the hierarchical control system presented as components in Chapters 5, 6, and 7 and presents the algorithmic operation of the system. Chapter 9 presents a simulation experiment that evaluates the hierarchical system against several

other ramp metering policies on a relatively small, but realistic, freeway management problem in the metropolitan Phoenix, AZ area. Presentation of performance variance information comparing metering methods has not been done before in freeway management literature. Finally, Chapter 10 summarizes the results of the research.

Chapter 2: Ramp metering literature review

Costs and benefits of ramp metering

Ramp metering is the most widely used form of freeway control [Yagar, 1989]. Ramp metering limits the rate at which vehicles enter the freeway system, thus potentially reducing the possibility of bottlenecks, shock wave propagation, and congestion. A wide range of benefits are available from the use of ramp metering [Arnold, 1987; Yagar, 1989; McShane and Roess, 1990]:

- (1) minimizing the total travel time of freeway users
- (2) efficient use of freeway capacity
- (3) discouraging routes with high societal costs
- (4) reducing the variance of corridor trip times
- (5) decreasing local freeway congestion and shock waves resulting from merging platoons
- (6) decreasing the accident rate in freeway weaving sections.

Another study indicates that ramp meters can efficiently reduce system travel time, although the savings are *network dependent* [Hellinga and Van Aerde, 1997]. Of course there are disadvantages and adverse effects of ramp meters:

- (1) encouraging longer trip distances on diversion routes
- (2) favoring through traffic over local traffic and short trips
- (3) modifying the evolved "status quo" of unobstructed freeway entry
- (4) increasing the overall operating cost of the control system
- (5) adversely affecting the surface street controller operation due to queue spillback and diversion to oversaturated locations

An operational study in the Denver area showed no statistically-significant improvement when a simple demand-capacity metering system was installed and evaluated [Lipp et al., 1991]. Ramp metering advantages may also be *strongly* dependent on the existence of good alternative routes, especially in the absence of effective queue management

strategies [Hellings and Van Aerde, 1997]. Nevertheless, most large metropolitan areas have some type of ramp metering installed or currently under installation, indicating that practitioners have been convinced that the benefits (fiscal, social, temporal) of ramp metering outweigh the costs of implementation, maintenance, and the adverse effects mentioned above.

Types of ramp metering algorithms

The vast array of ramp metering algorithms developed to date can be classified into one of three general categories;

- (1) *fixed-time* or time-of-day,
- (2) *traffic-responsive*, and
- (3) *hybrids* combining attributes of traffic-responsive and time-of-day algorithms.

Time-of-day metering algorithms

Time-of-day metering algorithms derive settings that apply during 10-30 minute intervals based on historical origin-destination flow rates and demand volumes for an entire commuting corridor or facility [Wattleworth and Berry, 1967; Messer, 1969; Yuan and Kreer, 1971; Wang, 1972; Wang and May, 1973; Chen et al., 1974; USDOT, 1976; Kahng et al., 1984].

The main drawback of fixed-time, time-of-day metering systems is the inability to handle non-recurrent incidents, accidents, special events, and fluctuations in traffic flow that may occur [Newman et al., 1970], since the actual demand may not be close to the demand used to derive the time-of-day metering rate. Recent studies have indicated that although time-of-day and day-to-day patterns exist, the variability of the actual flows from the historical average flows is significant enough to make some time-of-day settings ineffective [Rahka and Van Aerde, 1997].

Local traffic-responsive ramp metering algorithms

Traffic-responsive ramp metering algorithms measure variables such as speed, volume, and occupancy on the freeway and apply metering rates that keep the local freeway volume under capacity or at some desired set-point [Athans, 1969; Buhr et al., 1969;

Hardin, 1972; Estep, 1972; Pretty, 1972; HCM, 1985; Papageorgiou, 1989, 1991; Middelham and Smulders, 1991; Nihan, 1991; Nihan and Berg, 1992; Davis, 1993; Chang and Wu, 1994]. Other locally traffic-reactive ramp metering systems have been developed that merge vehicles into gaps in traffic [Drew et al., 1966; Wattleworth and Courage, 1968; Brewer et al., 1969] but such systems have not been widely implemented. Other types of traffic-reactive metering systems follow a pre-determined set of relationships between metering rates and traffic variable measurements. Examples of such systems are fuzzy and traditional rule-based expert systems and neural networks [Blumentritt et al., 1981; Rajan et al., 1986; Sasaki and Akiyama, 1987; Gray et al., 1990; Stephanedes et al., 1992; Zhang et al., 1994; Zhang and Ritchie, 1995; Papageorgiou et al., 1995].

The main drawbacks of using traffic-responsive ramp metering in a large-scale freeway are:

- (1) the absence of *coordination* between adjacent ramp meters, and
- (2) the absence of consideration of the *area-wide effects* of local changes to the metering rate.

Hybrid ramp metering control algorithms

Many hybrids and extensions of the basic traffic-responsive and time-of-day control methods have been developed. These extensions address the generally recognized issue that although day-to-day and time-of-day patterns exist, and can be exploited, their realization on a *specific* day and time may be significantly different from the assumed historical average pattern. Thus, hybrid ramp metering algorithms allow the system to follow the underlying trends, but still react to temporal and/or spatial flow irregularities.

The most straightforward extension of time-of-day methods is the use of a *rolling time horizon* and/or *periodic re-optimization* of the area-wide algorithm with new information [Messer, 1969; Drew et al., 1969; May, 1979; Papageorgiou, 1980, 1983; Kahng et al., 1984; Chang and Wu, 1994; Asakura, 1995]. However, the issue has been raised of how

long such re-optimization intervals should be. Most studies used *fixed* update intervals of 5-15 minutes.

The complexity of the large, multi-variable ramp metering problem has also been addressed by *decomposition* of the problem into smaller-scale descriptions of subsystems which are each optimized independently [Isaksen and Payne, 1973; Looze et al., 1978; Payne et al., 1979; Goldstein and Kumar, 1982; Papageorgiou, 1983; Kahng et al., 1984; Payne et al., 1985]. Many hybrids offer a combination of the rolling-horizon extension and spatial decomposition by establishing a *hierarchical* approach to the large-scale ramp metering problem, similar to the organizational structure shown in Figure 2-3 [Papageorgiou, 1983].

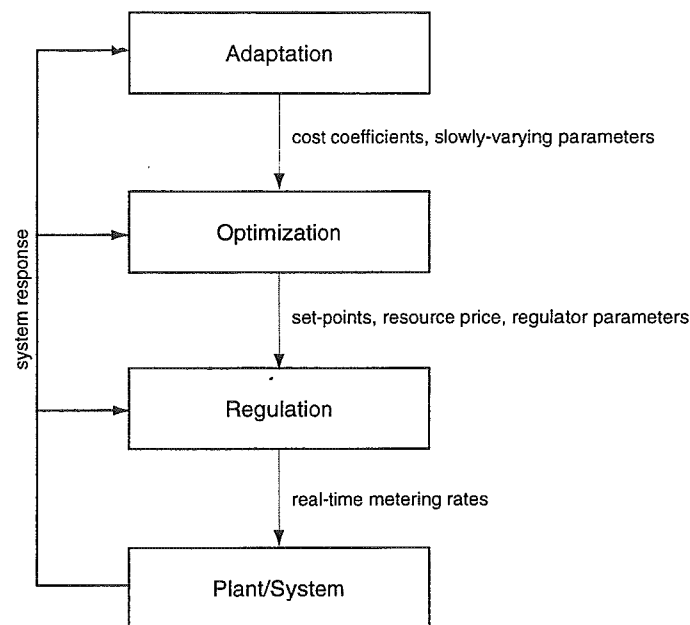


Figure 2- 1. Typical hierarchical control system structure

These systems combine locally-optimized traffic-responsive control with guidance from upper levels of the hierarchy regarding area-wide conditions, special events, and incidents [Drew et al., 1969; May, 1979; Papageorgiou, 1984; Payne et al., 1985]. Hybrid metering algorithms based on hierarchical control structures can also support the explicit consideration of optimization *modes* such as "normal flow", "congestion", and

"special-event" that are scheduled by the highest level(s) of the hierarchy [Papageorgiou, 1984; Pooran et al., 1994].

Integrated freeway/surface-street metering algorithms

Although there is a large body of work on freeway control algorithms, only exploratory work has been done to produce ramp metering solutions that integrate information from the surface-street system [Fan and Asmussen, 1990; Stephanedes and Chang, 1991, 1993; Pooran et al., 1992, 1994; Han and Reiss, 1994]. Some work has recently been proposed to develop freeway management solutions that derive *both* signal settings and metering rates in a commuting corridor of surface streets and freeway [Cremer et al., 1990; Chang et al., 1992; Papageorgiou, 1995; Zhang and Hobeika, 1997]. The failure to integrate the two systems has been due to the technological barriers that have restricted application of data-intensive ramp metering methods and the difficulty of modeling the two sub-systems together for optimization purposes [Van Aerde et al., 1987]. As such, no results of field implementation studies could be found in the literature.

However, as the technological barriers are being removed and real-time traffic information is becoming readily available, a new focus on improving the *system-wide* performance of the freeway and surface street network has emerged [Van Aerde and Yagar, 1988]. Critical data such as origin-destination (and/or route-proportional) matrices, time-varying demands, turning probabilities, and the like can be more reliably estimated on-line as the Intelligent transportation systems (ITS) "infrastructure" of communication networks and detection technology continues to be deployed.

Summary

For the past 35 years, much research has been done in the area of ramp metering control systems. Even simple metering systems installed in the field have been shown to be effective at improving freeway performance and having benefits that outweigh the installation and recurrent operating costs. However, metering systems sometimes have detrimental effects to the adjacent surface-streets when the ramp queue spills back into the interchange. Methods to address the spillback problem at the interface between the

freeway system and surface-street system have only recently been established in the research community and sparsely implemented in the field. Another drawback of ramp metering algorithms based on local traffic is the lack of consideration for the system-wide effects of the metering decisions and dis-proportionate queue growth rates [Benmohamed and Meerkov, 1994]. The remainder of this document describes a hierarchical freeway management system that builds on the successes of previous research in hybrid ramp metering algorithms and adds consideration of the important problem of queue management.

Chapter 3: Hierarchical ramp metering control system structure

Introduction

The ramp metering control system developed in this research is specifically designed to address the complicating factors of the freeway management problem. Recall from Chapter 1 that the freeway management problem is a difficult control and optimization problem because of these factors. Previous work in freeway ramp metering control systems has primarily focused on the complications caused by:

- (1) dynamic state changes,
- (2) stochasticity,
- (3) multi-dimensionality, and
- (4) partial observability

without consideration of *multiple objectives* or the *unpredictability* of the future traffic state.

The freeway control system developed in this research addresses all six of the complicating factors by establishing a hierarchical system of *layers* that

- (1) addresses *embedded* spatial and temporal descriptions of the ramp metering control problem,
- (2) considers concerns of *both* the freeway and surface-street systems in the optimization problem(s),
- (3) plans *pro-active* metering rates in real-time to respond to possible future traffic states, and
- (4) re-schedules optimizations based on the *stochastic fluctuations* of the demand processes.

Before detailing the characteristics of the hierarchical control model developed in this research, we review some concepts and previous research in hierarchical optimization.

Multi-level methods in hierarchical control

The hierarchical approach to system control has substantial fundamental research support, especially in the area of large-scale differential equation systems [Mahmoud,

1977; Sandell et al., 1978; Wilson, 1979; Papageorgiou and Schmidt, 1980; Bernassou and Titli, 1982; Papageorgiou, 1983]. Hierarchical control is particularly useful when the system being controlled has an appreciably large set of state variables, and/or and appreciably large set of control inputs. Large-scale control problems of this type are primarily difficult because of appreciable computation time required to solve for the “optimal” controls.

Such large-scale differential equation control problems are typically addressed by decomposing the problem using a multi-level approach. The multi-level approach creates a two-level optimization problem from a global optimization problem. Various methods have been proposed to solve the two-level optimal control problem including *interaction-prediction* and *interaction-balance* procedures [Sandell et al, 1978; Wilson, 1979; Papageorgiou and Mayr, 1982].

Multi-layer hierarchical control systems

A hierarchical control system can also describe a controller that solves the ramp metering problem at several embedded *layers* of aggregation. Thus, the multi-layer hierarchical approach typically indicates a structure where the targets, constraints, costs, and parameters of a given layer are communicated from a *higher-level* layer and the given layer communicates the targets, constraints, costs, and parameters to the *lower-level* layer(s) in the hierarchy [Mahmoud, 1977]. Few general theories exist to describe the effectiveness or expected performance of the multi-layer approach in system control since the definition and structure of such “layers” are problem-dependent [Sandell et al, 1978]. This approach has been implemented to address the freeway control problem with the layers being parameter estimation, incident detection, flow identification, and gap-acceptance ramp metering, respectively [Drew et al, 1969], although the gap-acceptance metering method has not found widespread acceptance. A later extension by Messer incorporated an LP-based area-wide coordination method at the “optimizing” layer of the hierarchy [Messer, 1971]. Modeling of the freeway control problem with a hierarchical, multi-layer approach has since persisted in the literature because of the natural way in which it addresses the complicating factors of the problem.

Set-point regulation methods

The hierarchical approach also applies to the development of *set-point* regulation control methods [Payne and Isaksen, 1973; Papageorgiou, 1983; Stephanedes and Chang, 1991]. A set-point regulation controller solves two separate optimization problems. One optimization problem (or problems) is solved to obtain the *set-point(s)* of the system. A second set of “optimization problems” are solved to obtain *control laws* that regulate the system state, under the influence of external disturbances, to operate at the set-points.

A third layer (*adaptation*) resides above the upper-layer optimization problem to modify the problem structure, parameters, and the like to the changing system conditions. The set-point regulation control method has been successfully applied in many areas of engineering such as chemical processing and aircraft control systems. Typically, because of the natural structure (i.e. geographic size, multi-dimensionality) of the system being controlled, the upper-layer control problem uses an *aggregate* model of the system to reduce processing requirements. Then at the lower-layer, the control problem is decoupled into *independent* subproblems that use a *more detailed* dynamic description of a *geographically smaller* portion of the problem given the assumptions of eqn. 3-3. Thus, because of the smaller size of the subproblems, more computational effort can be applied to solve each subproblem in real-time.

The MILOS hierarchical structure

This research addresses the system-wide ramp metering control problem by using a structured hierarchical framework hereafter referred to as the Multiobjective Integrated Large-Scale Optimized ramp control System (MILOS). This framework is based on the multi-layer approach to hierarchical process control using the set-point control method. Although neither the multi-layer approach to hierarchical control nor the set-point control method are contributions of this research, the hierarchical structure of the MILOS framework includes the following contributions:

- (a) consideration of *multiple objectives* in the optimization problem(s);

- (b) *integration* of information about the current conditions of the adjacent surface-street system;
- (c) *prediction* of possible future system states in the development of pro-active real-time metering rates; and
- (d) computability in real-time.

These primary characteristics of MILOS are driven by the structure of the real-world freeway/surface-street system. Thus, MILOS is composed of four hierarchically embedded, interactive subsystems:

- (1) locally reactive, predictive-cooperative real-time control,
- (2) area-wide coordination,
- (3) anomaly detection / optimization scheduling, and
- (4) subnetwork identification.

The structure of MILOS is a *pyramid* of modules that address smaller and smaller geographic areas of the large-scale ramp metering problems as one progresses lower in the hierarchy. The pyramid structure indicates that *one* optimization scheduler module schedules the solution of *several* area-wide coordination problems that in turn schedule the solution of several traffic-responsive real-time metering problems. This pyramid structure is illustrated in Figure 3-1, motivated by the structure of the RHODES hierarchical system for real-time *surface-street* traffic management and the RHODES-ITMS system developed at the University of Arizona [Head et al., 1992; Head and Mirchandani, 1993].

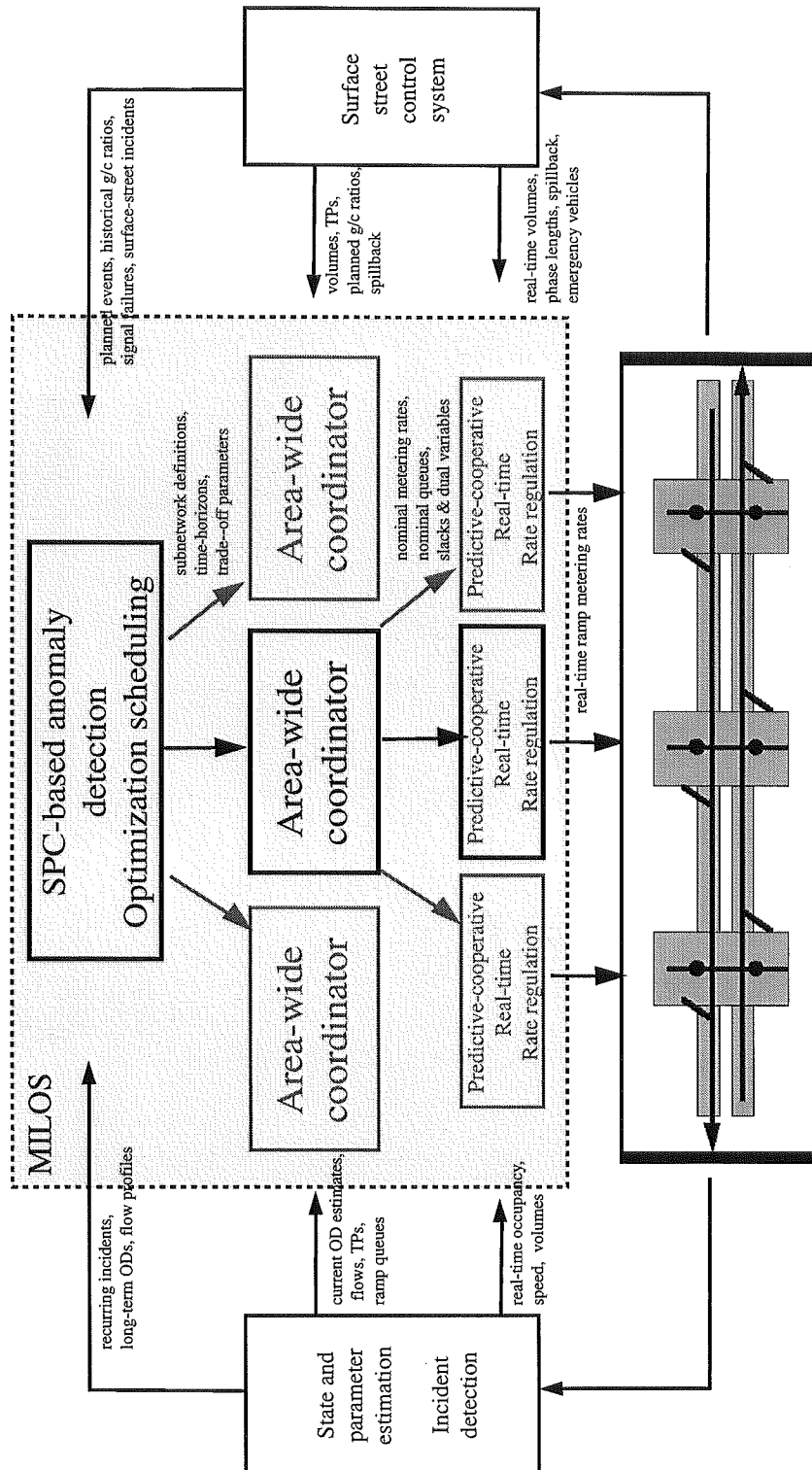


Figure 3- 1. The pyramid structure of the MILOS hierarchy

Modal decomposition of the MILOS hierarchy

MILOS can be considered to operate in several *modes*; strategic, tactical, and operational. At the highest level of the hierarchy, the *strategic* mode solves optimization problems with time horizons on the order of hours, days, and weeks, as well as responding to seasonal changes etc. The spatial influence of the strategic mode is the entire freeway and interchange network. The objectives of the strategic mode are to:

- (1) identify the “optimal” sub-network definitions that lower-level processors use to solve de-coupled optimization problems,
- (2) update parameters reflecting special events and long-term disturbances such as work zones,
- (3) update the slowly varying parameters in the system, and
- (4) determine the optimization time horizons for the lower-level problems.

The strategic mode is fulfilled by the *SPC anomaly detection* module and the *subnetwork identifier* module.

The *tactical* mode of the MILOS hierarchy solves optimization problems with time horizons of hours and minutes, using the subnetwork definitions and parameters passed from the strategic levels of the hierarchy. The spatial influence of a tactical-level module or optimization problem is “several ” (e.g. 5-15) adjacent ramp meters and the associated surface-street interchanges. The objectives of the *tactical* mode are to:

- (1) plan coordinated metering rates for recurrent congestion,
- (2) identify short-term flow fluctuations that require re-solution of the area-wide and real-time optimization problems,
- (3) react to changes in the relative congestion levels of the interchanges,
- (4) balance queue growth rates in a given geographic sub-network, and
- (5) respond to non-recurrent congestion generated by incidents.

The tactical mode is implemented by the *SPC-based anomaly detection* module and the *area-wide coordination* modules.

At the lowest level of the hierarchy, the *operational* mode solves optimization problems with time horizons of minutes, using the set-point metering rates and desired freeway

states provided by the tactical mode modules. The spatial influence of the operational level is a single ramp meter, a single interchange signal, and a small, *relatively predictable* subsection of the freeway. By “relatively predictable” it is meant that reasonable predictions for the next few minutes of flow can be made for this small section using a mathematical model. The objectives of the operational mode are to:

- (1) reduce ramp queue lengths when not detrimental to freeway conditions,
- (2) plan metering rates pro-actively based on prediction of possible future states,
- (3) react to short-term flow fluctuations that could cause freeway congestion, and
- (4) manage ramp queue spillback, if possible.

The operational mode is implemented by the *predictive-cooperative real-time control* modules that each solve optimization problems local to a single ramp meter.

Subnetwork identification

The majority of this research is focused on the area-wide optimization (tactical level) and traffic-responsive real-time control (operational level) algorithms. However, a role of the *strategic* mode in the MILOS framework is to identify the *problem boundaries* for the area-wide coordination and predictive-cooperative real-time control problems. Some research has been done to develop a method to determine boundaries for surface-street coordination problems [Moore and Jovanis, 1985], but little mention of such issues can be found in freeway control literature. For example, consider the large freeway network in Figure 3-2, where each node represents an interchange with a ramp meter (considering the unidirectional case only).

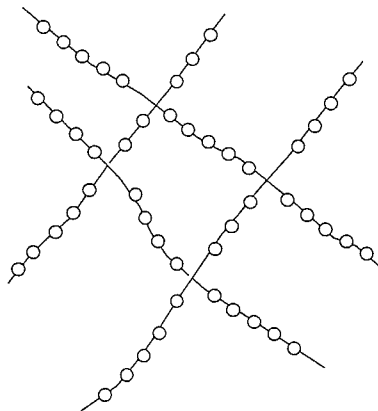


Figure 3- 2. Example freeway network

Consider an initial decomposition of the large-scale network problem into subsystems along each of the freeways, as shown by the boxes in Figure 3-3. Each of the subsystems would then be applied to a single area-wide coordination problem.

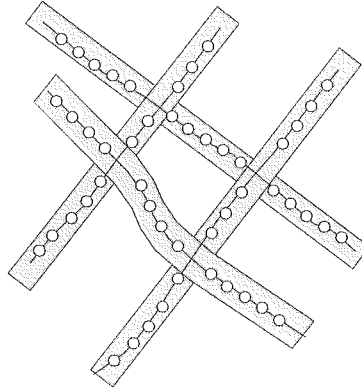


Figure 3- 3. Initial decomposition of freeway network

The question remains whether or not additional system performance could be gained by a different subnetwork structure, say for example, the structure indicated by Figure 3-4, where again each box indicates a separate area-wide coordination problem.

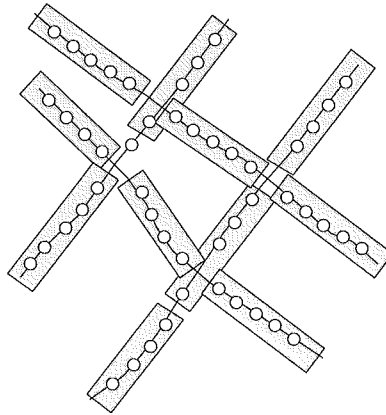


Figure 3- 4. Alternative decomposition of freeway network

No further development for the identification of subnetworks using analytical techniques was conducted in this research. Since many subnetwork definitions are easily, and “heuristically” pre-determined by the network topology, it was assumed that subnetworks for control is given. It should be noted that these subnetwork boundaries need not, and probably do not, coincide with political and jurisdictional boundaries [Fan and

Asmussen, 1990]. Thus, inter-agency cooperation would be necessary for cross-boundary coordination when a subnetwork crosses a boundary.

In this research, however, the subnetwork definitions at the area-wide coordination level will be taken as given by the traffic management decision-makers. It should be a topic of future work to develop the analytical *subnetwork identification* module of MILOS to advise freeway system managers of modifications to the subnetwork decomposition structure as current network conditions change.

SPC-based anomaly detection and optimization scheduling layer

The optimization procedures, at both the area-wide and real-time control layers, are continually re-evaluated using a *rolling-horizon* approach. Rolling-horizon approaches to traffic management have been proposed by many researchers in both surface-street and freeway control [May, 1979; Gartner, 1983; Chang et al, 1992; Head, et al., 1992; Sen and Head, 1997]. As the system evolves, the anomaly detector continually compares the observed freeway flows and ramp demands to the expected flows and demands. When a significant deviation from the expected state is detected, a new optimization run is scheduled immediately. The SPC anomaly detection module is based on the concept of control limits from the statistical process control (SPC) literature. This method is a completely novel approach to the “integrated” demand estimation and optimization scheduling problem and is discussed further in Chapter 7.

Area-wide coordination layer

The area-wide coordination layer provides the *tactical* decisionmaking of the MILOS hierarchy. The area-wide coordination level allocates medium-term (i.e. 10-20 minute) *target* or *nominal* ramp metering rates to maximize freeway throughput, balance ramp queue growth rates, and minimize queue spillback into the adjacent surface-street interchanges for a given subnetwork. The area-wide coordinator is based on a rolling-horizon implementation of a multi-criteria quadratic programming optimization problem. The area-wide coordinator interacts with the SPC anomaly detection module to identify over-capacity congestion conditions and to modify the optimization constraints and

criteria appropriately during incident conditions. The area-wide coordinator is sensitive to the needs of the adjacent surface streets by planning queue-growth rates according to the relative congestion level of each interchange. Several aspects of this formulation of the area-wide coordination problem are novel and discussed further in Chapter 5.

Predictive-cooperative real-time rate regulation layer

The predictive-cooperative real-time (PC-RT) rate regulation layer fulfills the *operational* mode of the MILOS hierarchy. The PC-RT optimization problems are based on a linearized description of the freeway state variables and are solved to minimize a linear measure of *additional* travel-time savings. This additional travel-time savings is above and beyond that due to the area-wide coordination solution by itself. At the *operational* layer, the system model is more detailed than at higher layers of aggregation [Papageorgiou, 1983; Payne et al., 1985; Fan and Asmussen, 1990]. Thus, linearization allows the PC-RT rate regulation module to plan, in real-time, several *pro-active* modifications to the nominal metering rates provided by the upper-layer area-wide coordination module based on predicted scenarios of possible ramp and freeway flows in the next few minutes. The scenario-based optimization structure of the PC-RT rate regulation module is a new treatment of the real-time ramp metering problem and its explicit connection to the solution of the area-wide coordination problem is completely new. These issues are discussed further in Chapter 6.

Integration of MILOS with necessary external systems

MILOS, as shown in Figure 3-1, is primarily an *optimization* system. Parameter estimation, especially turning-probability and route-proportional rate estimation, freeway detector data collection/filtering, incident detection and surface-street performance data all are taken as *inputs* to the MILOS hierarchy, and assumed to be “solved problems”. Of course, the successful implementation of an optimization routine such as MILOS is highly dependent upon the reliability and accuracy of the external algorithms and systems. In particular, MILOS requires real-time turning-probabilities, demand flows, green splits, and queue lengths from the interchanges control system. Such information

requires the availability of a sufficiently intelligent real-time signal controller and communications network.

Summary

A Multiobjective Integrated Large-Scale Optimized ramp control System (MILOS) is developed in this research. The framework is based on the multi-*layer* approach to hierarchical process control using the set-point/regulation paradigm. The MILOS framework is specifically structured to address the complicating characteristics of dynamic state changes, stochasticity, multi-dimensionality, partial observability, the existence of multiple objectives, and unpredictability that are inherent to the large-scale freeway control problem. In addition, MILOS considers the effects of freeway control decisions on the adjacent surface-street system at each level of the hierarchy. MILOS is composed of four hierarchically embedded, interactive subsystems:

- (1) area-wide coordination,
- (2) predictive-cooperative real-time control,
- (3) SPC-based anomaly detection and optimization scheduling, and
- (4) subnetwork identification,

based upon the decomposition of the large-scale control system into its *strategic*, *tactical*, and *operational* processing modes. In the next chapter, a popular and useful macroscopic flow model is discussed that is used (in Chapter 9) to evaluate the results of implementing the MILOS systems in a simulated freeway environment.

Chapter 4: Freeway Macrosimulator

Model construction

A macroscopic freeway traffic simulator based on the enhanced FREFLO [Payne, 1971, 1979; Rathi et al., 1985] and META models [Papageorgiou, 1984; Cremer, 1989] is used to evaluate and compare various ramp metering strategies developed in this research. The FREFLO macroscopic traffic simulator is based on partial differential equation (PDE) description of freeway traffic flow as a fluid of density $\rho(x,t)$ and speed $v(x,t)$ where x indicates spatial variation and t indicates temporal variation of the fluid's density and speed such that

$$\begin{aligned} \frac{\partial v}{\partial x} + \frac{\partial q}{\partial t} &= r - s \\ \frac{\partial v}{\partial t} &= v \frac{\partial v}{\partial x} - \frac{1}{T} \left[v - v_e(\rho) + v \frac{\partial \rho}{\partial x} \right] \end{aligned} \quad \text{Eqn. 4- 1}$$

where r is the ramp meter input rate, s is the off-ramp (or end of freeway) output rate $q(x,t)$ is the flow rate, $v_e(\rho)$ is the equilibrium speed-density relationship, and v is the anticipation coefficient [Lighthill and Witham, 1955; Richards, 1956; Michalopoulos et al., 1986, 1991, 1993]. This set of PDEs describes the conservation of vehicle flow through the freeway system and the dynamic relationship of speed and density. To evaluate $\rho(x,t)$ and $v(x,t)$ for various input rates r and exit rates s at each point along the freeway, the PDEs are discretized over space and time to obtain, using the simple Euler formula, the *difference* equation description of the system

$$\rho_j(k+1) = \rho_j(k) + \frac{T}{\Delta_j} (V_{IN,j}(k) - V_{OUT,j}(k) - s_j(k) + r_j(k)) \quad \text{Eqn. 4- 2}$$

$$\begin{aligned} v_j(k+1) &= v_j(k) + \frac{T}{\tau} (v_e(\rho_j(k)) - v_j(k)) + \frac{T}{\Delta_j} v_j(k) (v_{j-1}(k) - v_j(k)) \\ &\quad - \frac{v}{\tau} \frac{T}{\Delta_j} \left(\frac{n_{j+1}\rho_{j+1}(k) - n_j\rho_j(k)}{n_j\rho_j(k) + \kappa} \right) - \frac{\zeta T}{\Delta_j} \left(\frac{n_{j+1}r_{ON,j}(k)v_j(k)}{n_j\rho_j(k) + \kappa} \right) \end{aligned} \quad \text{Eqn. 4- 3}$$

$$\begin{aligned}
V_{IN,j}(k) &= \alpha \cdot V_{j-1}(\rho_{j-1}(k), v_{j-1}(k)) + (1 - \alpha) \cdot V_j(\rho_j(k), v_j(k)) \\
V_{OUT,j}(k) &= \alpha \cdot V_j(\rho_j(k), v_j(k)) + (1 - \alpha) \cdot V_{j+1}(\rho_{j+1}(k), v_{j+1}(k)) \\
V_{*,j}(\rho_j(k), v_j(k)) &= \rho_j(k) \cdot v_j(k)
\end{aligned} \tag{Eqn. 4- 4}$$

where $\rho_j(k)$ is the density, $v_j(k)$ is the mean speed, and $V_j(k)$ is the volume of vehicles in freeway section j at time k . Additional terms are added in (4-3) that are not represented in (4-1) for the speed PDE. $v_e(r_j(k))$ is an analytical speed-density characteristic such as

$$v_e(\rho_j(k)) = v_f \left(\left(1 - \frac{\rho_j(k)}{\rho_{MAX}} \right)^{l(3-2b_j)} \right)^m. \tag{Eqn. 4- 5}$$

The parameters v_f , ρ_{MAX} , l , and m of (4-5), as well as the other parameters of (4-2), (4-3), and (4-4), must be calibrated from field data and may vary over time and location. In particular, the time-interval T must be selected such that $T < \frac{\min[\Delta_j]}{v_f}$ to ensure that the state updates are frequent enough that flows do not “skip” sections.

For simplicity, we assume that the parameters do not vary from location to location during a given simulation. In addition, it is reasonable to assume that variations of the parameters v_f , ρ_{MAX} , l , and m of (4-5) are much slower than traffic flow dynamics and thus can be assumed as constant over a simulation period. Full description of the derivation of the remaining parameters in Table 4-1 can be found in [Payne, 1971; Cremer, 1989; Papageorgiou, 1989].

Symbol	Value
Δ_j	length of section j (km)
T	time interval duration (hr),
α	$\in [0,1]$, spatial discretization parameter
τ	time constant (km/hr), approximately the segment free-flow travel time
κ	constant (veh/km) to improve performance of eqn. 4-3 at low densities
v	anticipation coefficient (veh/km ²)
z	on-ramp friction coefficient
f	lane-drop friction coefficient
l	shaping parameter for speed-density characteristic
m	shaping parameter for speed-density characteristic
v_f	mean free-flow speed in section j (km/hr)
ρ_{MAX}	maximum density in a single traffic lane (veh/km)
ϕ	Influence factor of merging slowing effect $\in [0,1]$
n_j	number of lanes in section j
$r_j(k)$	on ramp rate in section j at time k (veh/hr)
$S_j(k)$	off ramp rate in section j at time k (veh/hr)
b_j	speed limit $\in [0,1]$ in section j (km/hr)
ζ	Influence factor of lane drop slowing effect $\in [0,1]$ (see eqn.4-6)

Table 4- 1. Parameters in macroscopic simulation equations

Equation 4-2 describes the evolution of the density $\rho_j(k)$ of each freeway segment j . Freeway sections which do not contain on-ramps have $r_j(k) = 0$ in (4-2) and sections without off-ramps have $s_j(k) = 0$ in (4-2). By convention, when both an off-ramp and on-

ramp are at an interchange, we define a single section that contains both the on and off-ramp.

Equation 4-3 describes the evolution of the speed $v_j(k)$ in link j over time. The four main terms of this evolution equation are included from both theoretical and empirical considerations [Payne, 1971]. The first term of (4-3) keeps the simulated speed from straying too far from the analytical speed-density relationship, and thus corrects for errors from the speed predicted by the analytical function $v_e(\rho_j(k))$. The second term is the "anticipation" term, indicating that speed in link j changes to reflect density changes *downstream* in link $j+1$ due to car-following behavior. The third term of (4-3) is the "convection" term that represents the effect of vehicle arrivals from upstream link $j-1$ on the speed in link j . The fourth term represents the slowing effect of merging vehicles from on-ramps and is > 0 iff $r_j > 0$ [Papageorgiou, 1989]. An additional term

$$-\frac{\phi T}{\Delta_j} \left(\frac{n_j - n_{j+1}}{n_j} \right) \left(\frac{\rho_j(k)}{\rho_{crit} n_j} \right) v_j(k)^2 \quad \text{Eqn. 4- 6}$$

was added to (4-3) by Papageorgiou to represent the slowing effect from a lane-drop when $n_{j+1} < n_j$. Previous research has indicated that this addition more accurately reflects this slowing effect than (4-3) without this term [Papageorgiou, 1989].

Equation 4-4 accounts for the spatial discretization of the flow model. The flow rate out of section j , $V_{OUT,j}(k)$, is expressed as a weighted sum of the flow rate $\rho_j(k)*v_j(k)$ from section j and section $j+1$, $\rho_{j+1}(k)*v_{j+1}(k)$, such that $\alpha \in [0,1]$. Similarly, the flow rate into section j $V_{IN,j}(k)$ is expressed as a weighted sum of the flow rates from section $j-1$, $\rho_{j-1}(k)*v_{j-1}(k)$, and j , $\rho_j(k)*v_j(k)$, to smooth the behavior of the model. Equation 4-4 does not, however, represent the general case where a segment j can have n feeder flows $V_{j,IN,1}(k), \dots, V_{j,IN,n}(k)$ and/or m receiver links $V_{j,OUT,1}(k), \dots, V_{j,OUT,m}(k)$. Such cases are straightforward additions to the single-source, single-receiver model by using weighted averages from the multiple sources/sinks for computing upstream and downstream state variables, ρ_{j+1} , v_{j+1} and ρ_{j-1} , v_{j-1} . respectively [Papageorgiou, 1984].

Modeling flow in heavy congestion

(4-2), (4-3), and (4-4) have been shown to accurately reflect freeway traffic when calibrated to a specific location during periods of moderate congestion [Papageorgiou, 1983, 1989; Cremer, 1989]. However, in situations of *heavy* congestion, the equations must be modified to reflect the fact that vehicles cannot continue to flow from section to section when a section is at the maximum density. The main reason for the continuing flow from section to section even though the congestion is high is that the dynamic equation for the section speed $v_j(k)$ does not accurately represent the breakdown in speeds when a section becomes congested. Previous researchers have addressed this by substituting the equilibrium speed-density relationship $v_e(\rho_j(k))$ for the dynamic speed-density relationship during periods of high congestion [Rathi et al., 1985]. We take a similar approach here, adding a threshold γ to (4-4) such that if the density exceeds the threshold, the density at the maximum capacity flow rate, the flow into or out of section j transitions to the theoretical volume-density relationship. Hence, during periods of high congestion and modify the structure of the density evolution equation while continuing to compute speeds using (4-3).

Essentially, we add a threshold γ to (4-2) such that the flow into or out of section j is equal to zero during sufficiently high congestion in the adjacent section. Hence,

$$\begin{aligned}
 V_{IN,j}(k) &= \alpha \cdot V_{j-1}(\rho_{j-1}(k), v_{j-1}(k)) + (1 - \alpha) \cdot V_j(\rho_j(k), v_j(k)) \\
 V_{OUT,j}(k) &= \alpha \cdot V_j(\rho_j(k), v_j(k)) + (1 - \alpha) \cdot V_{j+1}(\rho_{j+1}(k), v_{j+1}(k)) \\
 \text{if } \rho_j(k) \geq \gamma & \quad V_{*,j}(\rho_j(k), v_j(k)) = \rho_j(k) \cdot v_e(\rho_j(k)) \\
 \text{otherwise} & \quad V_{*,j}(\rho_j(k), v_j(k)) = \rho_j(k) \cdot v_j(k)
 \end{aligned}
 \tag{Eqn. 4-7}$$

Equation 4-7 is the modified form of (4-4). This condition helps to more accurately model the stagnancy of flow when density becomes overcritical and speeds are very low.

Special conditions must be added for the implementation of (4-7) for sections at the beginning and end of a freeway facility, so that if a congestion wave is passing upstream, it does not stagnate in the first segment and limit the input flow rate. Thus, since we do not have a measurement of $\rho_0(k)$, we use $\rho_1(k-1)$. Given that the congestion wave is passing upstream, the previous measurement of ρ_1 estimates the current density of ρ_0 , which is not available. This modification indicates that the source volume $V_0(k)$ must be reduced as the

congestion wave passes out of the system before it can return to the nominal value. Otherwise the density $\rho_i(k)$ in the first section will remain congested when the source volume $V_o(k)$ is reasonably large.

Dynamic modeling of ramp queues

In addition to the freeway state variables, it is also necessary to evaluate the queue lengths $q_i(k)$ at each ramp such that

$$\begin{aligned} q_i(k+1) &= q_i(k) + T(d_i(k) - r_j(k)) \\ q_i(k) &\geq 0 \end{aligned} \quad \text{Eqn. 4- 8}$$

where $r_i(k)$ is limited by

$$\begin{aligned} r_{MIN} &\leq r_j(k) \leq \hat{r}_{MAX} \\ \hat{r}_{MAX} &= \begin{cases} \min(d_i(k), r_{MAX}) & \text{if } q_i(k) = 0 \\ r_{MAX} & \text{otherwise} \end{cases} \end{aligned} \quad \text{Eqn. 4- 9}$$

where $r_{i,MIN}$ and r_{MAX} are given minimum and maximum ramp metering rates, respectively. Since $r_i(k)$ can be set higher than the demand rate when a queue is present (as high as the saturation flow rate), we must restrict $q_i(k)$ to be non-negative since we cannot have negative queues. At this level of modeling, we do not consider driver behavior in the metering rate limitations. Thus, when a given metering rate is specified, (e.g. 456 veh/hr) it is assumed that drivers can implement this rate precisely (e.g. not 445 veh/hr or 500 veh/hr). Recent results have indicated that this is usually not true in the real system, especially at high metering rates when reaction time may be very close to the allocated green-time of each metering signal [Banks, 1992; Decker, 1997].

To reflect the fact that vehicles are slowed, and possibly stopped, when entering and exiting a freeway segment that is highly congested, we can add a soft-limiter

$$\xi(\rho_j(k)) = \frac{e^{\left(\frac{\rho_c - \rho_j(k)}{\rho_{MAX} - \rho_c}\right)}}{1 + e^{\left(\frac{\rho_c - \rho_j(k)}{\rho_{MAX} - \rho_c}\right)}} \quad \text{Eqn. 4- 10}$$

to the on-ramp rate $r_i(k)$ and off-ramp rate $s_i(k)$. The shape of (4-10) is a sharp, but continuous, transition about the critical point ρ_c as illustrated in Figure 4-1.

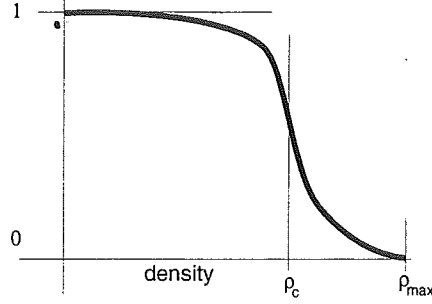


Figure 4- 1. Typical shape of soft-limiter function

The limiting effect is also applied to the off-ramp rates, since, under congested conditions, the off-ramp is also blocked after some critical density ρ_c is exceeded. Previous models without such blockage terms would underestimate the clearance time required for congestion to dissipate. Hence, (4-2) is further modified as

$$\rho_j(k+1) = \rho_j(k) + \frac{T}{\Delta_j} \left(V_{IN,j}(k) - V_{OUT,j}(k) + \xi(\rho_j(k)) (r_j(k) - s_j(k)) \right) \quad \text{Eqn. 4- 11}$$

and (4-8) as

$$\begin{aligned} q_i(k+1) &= q_i(k) + T \left(d_i(k) - \xi(\rho_j(k)) \cdot r_j(k) \right) \\ q_i(k) &\geq 0 \end{aligned} \quad \text{Eqn. 4- 12}$$

to reflect the fact that a queue can develop at the ramp, even if the ramp metering rate is higher than the demand rate, when freeway congestion blocks vehicles merging from the ramp. ρ_c and ρ_{MAX} must be carefully chosen (calibrated) in equation to reflect realistic effects of queue-growth and restricted flow in the presence of freeway congestion. Nevertheless, their addition to the macro-simulation model is to more accurately represent congested conditions for comparison of various metering strategies in the evaluation experiment of Chapter 9. Some preliminary evidence of the positive effects of these additions are also shown in the benchmark test example of the area-wide coordination problem in Chapter 5.

Summary of macroscopic model

The macroscopic description of freeway traffic and ramp queues used to evaluate ramp metering strategies in this research project is a system of nonlinear difference equations based on the fluid-flow model given by

$$\rho_j(k+1) = \rho_j(k) + \frac{T}{\Delta_j} \left(V_{IN,j}(k) - V_{OUT,j}(k) + \xi(\rho_j(k)) (r_j(k) - s_j(k)) \right)$$

$$\begin{aligned}
v_j(k+1) &= v_j(k) + \frac{T}{\tau} (v_e(\rho_j(k)) - v_j(k)) + \frac{T}{\Delta_j} v_j(k) (v_{j-1}(k) - v_j(k)) \\
&\quad - \frac{\nu}{\tau} \frac{T}{\Delta_j} \left(\frac{n_{j+1}\rho_{j+1}(k) - n_j\rho_j(k)}{n_j\rho_j(k) + \kappa} \right) - \frac{\zeta T}{\Delta_j} \left(\frac{n_{j+1}r_{ON,j}(k)v_j(k)}{n_j\rho_j(k) + \kappa} \right) \\
V_{IN,j}(k) &= \alpha \cdot V_{j-1}(\rho_{j-1}(k), v_{j-1}(k)) + (1 - \alpha) \cdot V_j(\rho_j(k), v_j(k)) \\
V_{OUT,j}(k) &= \alpha \cdot V_j(\rho_j(k), v_j(k)) + (1 - \alpha) \cdot V_{j+1}(\rho_{j+1}(k), v_{j+1}(k)) \\
\text{if } \rho_j(k) \geq \gamma \quad &V_{*,j}(\rho_j(k), v_j(k)) = \rho_j(k) \cdot v_e(\rho_j(k)) \\
\text{otherwise} \quad &V_{*,j}(\rho_j(k), v_j(k)) = \rho_j(k) \cdot v_j(k) \\
q_i(k+1) &= q_i(k) + T(d_i(k) - \xi(\rho_j(k)) \cdot r_j(k)) \\
q_i(k) &\geq 0 \\
r_{MIN} \leq r_j(k) &\leq \hat{r}_{MAX} \\
\hat{r}_{MAX} &= \begin{cases} \min(d_i(k), r_{MAX}) & \text{if } q_i(k) = 0 \\ r_{MAX} & \text{otherwise} \end{cases} \\
\xi(\rho_j(k)) &= \frac{e^{\left(\frac{\rho_c - \rho_j(k)}{\rho_{MAX} - \rho_c}\right)}}{1 + e^{\left(\frac{\rho_c - \rho_j(k)}{\rho_{MAX} - \rho_c}\right)}} \\
v_e(\rho_j(k)) &= v_f \left(\left(1 - \frac{\rho_j(k)}{\rho_{MAX}} \right)^{l(3-2b_j)} \right)^m
\end{aligned}$$

This model has the parameters $[\rho_c, \varphi, \rho_{crit}, \kappa, \alpha, \tau, \zeta, l, m, \rho_{max}, v_p, T]$, the specific geometric and travel-behavior details $[d_{IN,j}(k), s_{OFF,j}(k), \Delta_j, n_p, b_j]$, and the control variables $r_j(k)$. The parameters *must* be calibrated precisely to the specific location characteristics and driving population of the intended application area to obtain reasonable real-world flow behavior from the model. Inappropriate choices for the model parameters can easily lead to “unstable” model performance and inconclusive results. This model has been modified slightly from previous instances of the model to reflect ; (a) the condition where queues are built at the freeway on-ramps when the density in a section is too high to allow the current on-ramp flow rate, and (b) the condition that off-ramp rates are also reduced when the density becomes large because vehicles cannot physically move to the off-ramp to exit the freeway system when the speed is near zero.

Stochastic effects and diversion behavior

Note that this macroscopic description of freeway flow is a deterministic model. Empirical data for real-world freeways indicates that the system does *not* evolve in a completely deterministic manner, but is highly affected by stochastic disturbances and flow fluctuations. We address this issue by treating the input streams $d_i(k)$ and the initial upstream freeway input(s) $V_o(k)$ as random variables, but leaving the evolution equations deterministic, as opposed to previous approaches based on adding acceleration noise to the dynamic speed equations [Weits, 1988]. It is shown empirically in Chapter 9 that considering the inputs $d_i(k)$ and $V_o(k)$ as random variables significantly improves the match of the macroscopic model to a stochastic, *microscopic* model of a study area in Phoenix, AZ. This microscopic simulation model, CORSIM, simulates travel of individual vehicles in one-second increments and is well accepted for extensive simulation testing and evaluation.

This macroscopic flow model also does not *explicitly* simulate diversion behavior or route modification, but this can be “easily” added by modifying/updating the route-proportional matrix (which determines the off-ramp rates $s_j(k)$ and demands $d_i(k)$) due to the current conditions. Of course, as has been indicated in previous work, estimation and re-estimation of route-proportional matrices and diversion rates is very difficult and is a subject of much research [Cremer and Keller, 1987; Madanat et al., 1995; Ashok and Ben-Akiva, 1993; Ding et al., 1996]. However, the ultimate success of a ramp metering control system such as MILOS in real-world freeway systems is highly dependent upon accurate and reliable turning-proportions and/or route-proportional matrix estimation. This research will assume that route-proportional matrices are given. In Chapter 7 we present a method that *could* be used to detect changes to the route proportions and/or turning-probabilities, but we do not further explore this possibility.

Summary

A macroscopic freeway traffic simulator based on the enhanced FREFLO and META models was developed for evaluation of various ramp metering strategies developed in this research project. The performance of the model was improved to represent highly congested conditions, especially the simulation of ramp queues. A term was added to the flow equations that represents the inability of vehicles to enter the freeway from the ramp when the freeway is so congested that no merging maneuver can occur. Specific simulation results showing the effects of the modeling enhancements are presented in

Chapter 9. The next three chapters develop the area-wide coordination, locally-reactive real-time optimization, and SPC-based anomaly detection and optimization scheduling layers of the MILOS hierarchical control structure.

Chapter 5: Area-wide Coordination Problem

Introduction

The area-wide coordination layer provides primarily the *tactical* decision-making of the MILOS hierarchy by providing *target* ramp metering rates based on area-wide conditions and aggregate traffic flows in each segment. The area-wide coordinator is based on two rate coordination problem formulations from the literature

- (1) Yuan and Kreer's queue-balancing problem [Yuan and Kreer, 1971] and
- (2) Wattleworth and Berry's throughput maximization problem [Wattleworth and Berry, 1968].

These formulations are significantly modified in the model presented here. By using a multi-criterion objective function, we combine the two conflicting objectives to address both total system performance and user-specific performance benefits. The objective function also includes

- (a) consideration of the specific differences in interchange congestion,
- (b) physical capacity along the corridor, and
- (c) agency/system-operator preference for incorporating queue-growth considerations [Fan and Asmussen, 1990].

Mathematical description of the area-wide coordination problem

Consider a unidirectional freeway with N on-ramps and M off-ramps where the demand (veh/hr) d_i at each on-ramp i , $i=1...N$ is provided by either

- (1) the physical beginning of the freeway facility,
- (2) a freeway-freeway connector, or
- (3) a surface-street interchange.

We assume that demand $d_0=V_0$ provided at the beginning of the freeway cannot be controlled via ramp metering. The only freeway controls available are ramp metering rates (veh/hr) r_i $i=2...N$. By convention we assume that $r_1 = d_0$, (i.e. the freeway input or the first "ramp" is uncontrollable. Speed limits are assumed fixed in each freeway section, but need not be equal everywhere. Speed advisories, such as those that could be provided by variable message signs (VMS), are not considered in this algorithm.

The vehicular flows x_j in each freeway link j are determined by evaluating the route-proportional flows from each on-ramp to each off-ramp, such that

$$\sum_{i=1}^j A_{i,j} r_i = x_j \quad \forall j \quad \text{Eqn. 5- 1}$$

where $A_{i,j}$ values have the special structure

$$\begin{aligned} 0 &\leq A_{i,j} \leq 1 \\ A_{i,j} &= 0 \quad i < j \\ A_{i,j} &\leq A_{i,j-1} \quad i > j \end{aligned} \quad \text{Eqn. 5- 2}$$

The matrix $A = \{A_{i,j}\}$ describes the proportion of the flow entering at ramp i that continues through link j en route to its destination, it will be referred to as the route-proportional matrix. The matrix A is assumed to be constant and known over the control period horizon, T . In a *steady-state* input-output description of the freeway system such that $x_j(k) = \bar{x}_j \quad \forall k \leq T$, it is assumed that all demand entering at ramp i bound for off-ramp j will exit at off-ramp j during the time horizon T . Thus, we need only be concerned with the physical limit of freeway capacity

$$\sum_{i=1}^j A_{i,j} r_i \leq CAP_j \quad \forall j \quad \text{Eqn. 5- 3}$$

in each segment. The physical limit CAP_j is derived for each segment from the volume-density curve specific to that segment. The volume-density curve can be empirically derived (i.e. curve-fit) from observations or computed from the saturation flow rates, number of lanes, merge area restrictions, and other factors as detailed in established procedures [McShane and Roess, 1990]. Given the concerns noted in Chapter 1, it may be advantageous to set a capacity CAP_j for each link j in the optimization model that is slightly less than the *critical* maximum volume to ensure stable flow. As will be shown in Chapter 9, it is difficult to maintain flows at the critical value CAP_j without beginning a backward-traveling congestion wave, confirming the difficulties of the freeway management problem as presented in Chapter 1.

An additional necessary set of constraints for the area-wide coordination problem is a limitation on the minimum and maximum ramp metering rate, such that

$$r_{i,MIN} \leq r_i \leq r_{i,MAX} \quad , \quad i = 1 \dots N \quad \text{Eqn. 5- 4}$$

where $r_{i,MAX} = \min(d_i, s_i)$. Here, s_i is the saturation flow rate of the ramp i (a ramp could have more than one lane) and $r_{i,MIN}$ is the slowest rate acceptable to drivers, such as two vehicles per minute (120 veh/hr). The rate $r_{i,MIN}$ could be as low as zero if the ramp was allowed to be and/or capable of being fully closed.

In this optimization formulation, metering/closure of a ramp only creates a queue at the ramp and does *not* result in driver diversion. Diversion rates would be computed by an external processor (not discussed in this research) that updates the route-proportional matrix and demands given the control decisions, e.g. [Cremer and Keller, 1987; Ashok and Ben-Akiva, 1993; Madanat et al., 1995; Ding et al., 1997]. Chapter 7 discusses a system identification procedure based on statistical process control which could be used to aid diversion modeling by detecting changes to the flows x_j in each link deviating from their assumed nominal values \bar{x}_j .

Derivation of the objective function

Given the constraints detailed above, a popular objective function is to maximize the total inputs to the freeway

$$\max_{r \in R} \sum_{i=1}^N r_i. \quad \text{Eqn. 5- 5}$$

This objective is derived from minimizing the total travel time in the freeway system, which is a typical operational goal of freeway control [Wattleworth and Berry, 1968; Messer, 1971; May, 1979; Papageorgiou, 1983]. Using this objective, the current freeway conditions $\rho_j(k)$, and on-ramp demands $d_i(k)$ must be continually monitored and compared to the assumed steady-state values $\bar{\rho}_j$ and d_i . When the values of $\rho_j(k)$ and $d_i(k)$ drift outside of a reasonable upper or lower bound, the problem must be redefined and re-optimized [Messer, 1971; May, 1979] as developed in Chapter 7.

Consideration of queue storage limits

The classical linear programming rate coordination problem is described by (5-3), (5-4), and (5-5) [Wattleworth and Berry, 1968]. This formulation does *not* consider the formation of ramp queues $q_i(k+1) = q_i(k) + \Delta T(d_i(k) - r_i(k))$ at each on-ramp due to the

application of metering rates less than the offered demand. Thus, to reflect the physical limitations of ramp queuing areas, an additional set of constraints must be added such that

$$(d_i - r_i)T \leq Q_i \quad \forall i \quad \text{Eqn. 5- 6}$$

such that Q_i is the physical limit on the number of vehicles that can be stored on the ramp without causing spillback into the interchange (assuming some average vehicle length) and T is the optimization time horizon. Q_i would be based on the length of the ramp storage area and the average vehicle length. In operational practice, one may want to keep the capacity limitation Q_i slightly less than the physical limit of storage to provide an additional cushion for unexpected surges in demand. The constraints (5-6) limit the rate at which the queue is allowed to grow (based on a constant arrival rate) and *fill to capacity* during the optimization horizon T . Vehicles queued at the ramp at the beginning of the optimization period are included in the offered demand d_i such that

$$d_i = d_{ext} + \frac{q_i(0)}{T} \quad \forall i \quad \text{Eqn. 5- 7}$$

converting the queued vehicles $q_i(0)$ into a flow rate (veh/hr) by assuming that all of the vehicles queued “demand” to be discharged during the time horizon T .

Inclusion of constraints (5-6) would indicate that, in the absence of re-optimization during the time horizon T , at $t = t_0 + T$, several queues *may* be filled to capacity. This would require, at least for a short time, $r = r_{i,MAX}$ (saturation flow rate) to clear the queue and to create ramp capacity. This clearing at the maximum rate could have significant detrimental effects to freeway conditions as the vehicles attempt to merge into traffic as a platoon. We address the problem of alternately filling queues to capacity and dissipating them at the saturation flow rate in two ways. First, we implement a *rolling-horizon* solution to the area-wide coordination problem (5-3, 5-4, 5-5, 5-6) with frequent estimates of the current constant demand d_i and queue lengths q_i . Thus, as the unfilled queue storage capacity begins to decrease as q_i approaches Q_i , the demand rate at the ramp increases and it becomes more likely that queue dissipation will occur. Second, by modifying the nominal rate $r_{i,N}$ in real-time such that $r_i(k) = r_{i,N} + \Delta r_i(k)$ by solving a predictive-cooperative optimization problem at each ramp for $\Delta r_i(k)$, we can take advantage of the opportunities to *dissipate queues* when

- (a) the demand rate to the ramp meter is lower than expected and/or
- (b) the freeway conditions are lighter than expected.

More details of the predictive-cooperative real-time optimization subproblems solved at each ramp are provided in Chapter 6.

Development of a multi-criteria objective function

Inclusion of constraints (5-6) into the linear programming problem formulation provides queue-growth management due to physical limitations of each ramp, but does *not* control queue-growth according to the prevailing *congestion levels* at each interchange. Such constraints (5-6) also do *not* guarantee that equitable decisions will be made as to where to hold vehicles in queues to provide freeway congestion relief. A quadratic optimization criterion

$$\min_{r \in R} \sum_{i=1}^N (d_i - r_i)^2 \quad \text{Eqn. 5- 8}$$

was proposed by Yuan and Kreer to address the need to balance ramp queues at each ramp, such that $q_1 \cong q_2 \cong \dots \cong q_{n-1} \cong q_n$ (rather than hold many vehicles at some ramps and none at others, such that $q_1 \gg 0$, $q_2 \gg 0$, $q_3 = \dots = q_{n-1} = q_n = 0$ [Yuan and Kreer, 1971] which is a typical result of linear programming formulations such as (5-3), (5-4), and (5-5) where the objective results in optimal solutions at the extrema of the feasible region). We can thus use a combination of the two objectives (5-5) and (5-8) to obtain a compromise solution that addresses *both* freeway throughput and ramp queue management.

Thus, we would like to simultaneously minimize freeway total travel time (by maximizing on-ramp flow in steady-state) and balance ramp queues throughout the corridor. It would be imprudent to simply add the cost functions (5-5) and (5-8) together because the units are not the same (i.e. (veh/hr) and (veh/hr)², respectively). As such, we use a simple technique to combine objectives with differing units by dividing each objective by the "ideal" cost and adding the dimensionless quantities. However, in (5-8) the optimal cost is zero when $r_i^* = d_i$ $i=1\dots N$, and thus we cannot divide by the ideal cost solution to obtain a dimensionless objective for (5-8).

Thus, we modify objective (5-8) from *minimizing* the distance from the ideal point $r_i^* = d_i$ $i=1\dots N$, to *maximizing* the distance from the *anti-ideal* point. The anti-ideal point is

the (also usually infeasible) solution $r_i^* = r_{i,MIN}$ $i=1...N$ which creates the longest possible queues at each ramp, and thus, the worst possible value for (5-8). Hence our single objective now combines the distance for (5-5) from the ideal point $r_i^* = d_i$, $i=1...N$ and the distance for (5-8) from the anti-ideal point $r_i^* = r_{i,MIN}$ $i=1...N$ resulting in a compromise objective function

$$\max_{r \in R} \sum_{i=1}^N r_i + \left(\sum_{i=1}^N d_i \right) \left(1 - \frac{\sum_{i=1}^N (d_i - r_i)^2}{\sum_{i=1}^N (d_i - r_{i,MIN})^2} \right) \quad \text{Eqn. 5- 9}$$

Even though the terms from (5-5) and (5-8) are now in equivalent units, the *relative* difference in the size of the two cost components for typical feasible choices of r_i will *still* influence in the importance attributed to each objective. We can provide decision-makers with a preference between the two components by including a weighting factor β such that (5-9) becomes

$$\max_{r \in R} \sum_{i=1}^N r_i + \beta \left(\sum_{i=1}^N d_i \right) \left(1 - \frac{\sum_{i=1}^N (d_i - r_i)^2}{\sum_{i=1}^N (d_i - r_{i,MIN})^2} \right). \quad \text{Eqn. 5- 10}$$

Setting β large will increase the importance of balancing ramp queues and setting β small will decrease the importance on balancing queues and increase the importance of maximizing freeway throughput.

Setting costs according to interchange congestion level

Although the objective (5-14) includes considerations for queue growth, the mechanism to distinguish queue growth at one ramp over another is only the storage limitations in (5-6) and the freeway conditions surrounding each ramp location. To reflect the current congestion conditions at each interchange, we weight each of the components of objective (5-10) with a weighting factor c_i such that

$$c_i = \frac{\sum_{m=1}^M \frac{V_{m,i}}{C_{m,i}}}{\max_i (c_i)} \quad \forall i \quad \text{Eqn. 5- 11}$$

where $C_{m,i}$ is the capacity of phase m at interchange i and $V_{m,i}$ is the offered volume for phase m at interchange i . Thus, the weighting factors c_i reflect the relative importance of vehicle storage on one ramp versus another according to the possible impacts on the surface streets if spillback should occur. This result was developed independently in this research, but found to have been proposed previously for the queue balancing objective (5-8) [Fan and Asmussen, 1990]. As the congestion levels change at each interchange, the cost coefficients c_i are updated to reflect the most current conditions. Such updates are enacted at least as often as the area-wide coordination optimization problem is re-solved. However, the scheme for updating c_i estimates is not in the scope of this research.

The weighting factors c_i could also be set by decision-makers/system-operators based on other considerations such as

- (a) ad-hoc values set to discourage or encourage long-term flow changes at certain interchanges,
- (b) average delay at each interchange,
- (c) surface-street incident conditions, as well as
- (d) virtually any other performance measure computable in real-time, and as suggested by agency preference [Powell, 1997].

In this research, we restrict derivation of the c_i cost coefficients to the congestion level as reflected in the time-varying V/C ratio.

Integration of surface-street flows in ramp demands

Use of the weighting factors c_i derived in real-time as the area-wide coordination optimization problem is re-solved is one way in which the area-wide coordination optimization problem incorporates the interchange conditions. In addition, the area-wide ramp metering coordination optimization problem incorporates interchange conditions in its decision-making by building the ramp demands d_i from the surface street flows using

$$d_i = p_{R,NB}d_{NB} + p_{L,SB}(1 - p_{R,SB})d_{SB} + p_{T,EB}d_{EB} + \frac{q_i(0)}{T} \quad \forall i \quad \text{Eqn. 5- 12}$$

where $q_i(0)$ is the queue length at the ramp when the optimization begins, $p_{R,NB}$, $p_{L,SB}$, and $p_{T,EB}$ are the current probabilities of turning right, left, and through, respectively at each of the approaches to the interchange feeding ramp i , and d_{NB} , d_{SB} , and d_{EB} are the demands on the northbound, southbound, and eastbound approaches to interchange i , respectively.

These definitions assume an Eastbound freeway simply for demonstration purposes, as illustrated in Figure 5-1.

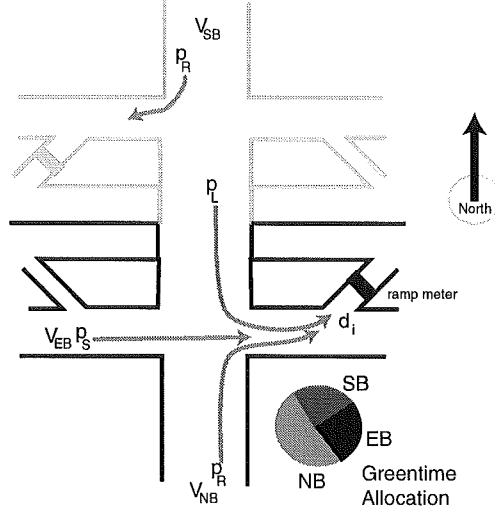


Figure 5- 1. Ramp meter demand sources

Using the turning probabilities $p_{R,NB}$, $p_{L,SB}$, and $p_{T,EB}$ and surface-street demands d_{NB} , d_{SB} , and d_{EB} to build the ramp demand d_i makes the assumption that these quantities are reasonably constant (and available) over the optimization time horizon.

It should be noted that although the method(s) chosen to derive the turning probabilities $p_{R,NB}$, $p_{L,SB}$, and $p_{T,EB}$ and estimate the real-time demands d_{NB} , d_{SB} , and d_{EB} will play a significant role in the validity and applicability of the resulting area-wide rate coordination, the coordination algorithm is *independent* of the method used. Slight inaccuracies in turning probability and demand rate estimation should be mitigated because the area-wide rates r_i are modified in real-time by the lower-layer optimization sub-problems. In the event that surface street demand and/or turning probability estimates are *significantly* different from current real-time conditions, the SPC-based anomaly detection module will identify this and begin a new iteration of the area-wide coordination optimization problem with appropriately modified parameters d_{NB} , d_{SB} , and d_{EB} and $p_{R,NB}$, $p_{L,SB}$, and $p_{T,EB}$. The operation of the SPC module is detailed in Chapter 7.

Quadratic objective function summary

To simplify notation of the objective function of the area-wide coordination optimization problem, (5-10) with the additions of (5-11) can be expressed as

$$\max_{r \in R} \sum_{i=1}^N r_i + \beta \gamma \left(\sum_{i=1}^N c_i d_i - c_i (d_i - r_i)^2 \right) \quad \text{Eqn. 5- 13}$$

where

$$\gamma = \frac{\left(\sum_{i=1}^N d_i \right)}{\sum_{i=1}^N c_i (d_i - r_{i,MIN})^2} \quad \text{Eqn. 5- 14}$$

Expanding the square and neglecting terms that do not contain the decision variables r_i we obtain the optimization problem

$$\max_{r \in R} J = \sum_{i=1}^N (1 + 2\beta\gamma c_i d_i) r_i - \beta\gamma c_i r_i^2 \quad \text{Eqn. 5- 15}$$

subject to the constraints (5-3), (5-4), and (5-6) where c_i is specified as in (5-11), d_i is specified as in (5-12), and all other parameters A , β , T , Q_i , $r_{i,MIN}$, $r_{i,MAX}$ and $q_i(0)$ are specified from external data. (5-15) is a quadratic objective and (5-3), (5-4) and (5-6) are linear constraints and thus the solution has a unique optimum when a feasible solution exists.

Resolving infeasibility

It is possible, however, that the formulation posed in (5-3), (5-4), (5-6), and (5-15) does *not* have a feasible solution. For example, an accident on a freeway link could reduce the capacity considerably in that section, requiring many more vehicles to be metered at upstream ramps than could be stored in the available ramp queues. In such a case, constraints from (5-3), (5-4), or (5-6) must be relaxed to render the problem feasible. Constraints (5-6) are the best candidate for relaxation, since we cannot increase the physical carrying capacity of a freeway section in (5-3) or change the maximum (minimum) possible metering rate, limited by the saturation flow rate (zero), in (5-4). Thus, we must opt to allow spillback for a short time into the interchanges by increasing the queue storage capacity in constraints (5-6) to accommodate the overflow in a system-equitable manner.

Let $z_i \geq 0$ $i = 1 \dots N$ be the extra capacity allocated at each ramp queue i to accommodate the flow at that ramp. In the same way that the queue storage is balanced according to the

interchange congestion cost c_i in a feasible problem, consider an allocation of the queue *overflow* in a similar manner. Thus, the constraints (5-6) are modified such that

$$(d_i - r_i)T - z_i \leq Q_i \quad \forall i. \quad \text{Eqn. 5- 16}$$

A penalty term is added to the objective function (5-15) incorporating the cost of allowing queue i to extend beyond its capacity Q_i over the time horizon T , such that the objective function now becomes

$$\max_{r \in R} J = \sum_{i=1}^N (1 + 2\beta\gamma c_i d_i) r_i - \beta\gamma c_i r_i^2 - \beta_2 \gamma c_i z_i^2 \quad \text{Eqn. 5- 17}$$

where β_2 is an appropriately chosen scaling constant. In particular, β_2 should be specified large enough, say $\beta_2 = 100 \beta$, to induce $z_i \cong 0$ for all solutions that are feasible *without* the inclusion of the additional capacity variables z_i , $i=1 \dots N$. Choosing a "small" value of β_2 can result in a solution where some ramps are specified to spill-back, and others are allowed to flow unconstrained. This is much like the solution resulting from using an LP method, but with the c_i terms the congestion at each interchange is considered. It should be noted here that the objective function (5-17) has no physical meaning with the introduction of the penalty term $\beta_2 \gamma c_i z_i^2$ and is likely a negative quantity in the overcapacity situation. However, each of the components of (5-17) is suitably derived to benefit both freeway and surface-street system operation.

Area-wide coordination problem summary

The quadratic optimization problem is summarized as

$$\max_{r \in R} \sum_{i=1}^N (1 + 2\beta\gamma c_i d_i) r_i - \beta\gamma c_i r_i^2 - \beta_2 \gamma c_i z_i^2 \quad \text{Eqn. 5- 18}$$

subject to

$$\sum_{i=1}^j A_{i,j} r_j \leq CAP_j \quad \forall j$$

$$(d_i - r_i)T - z_i \leq Q_i \quad \forall i$$

$$r_{MIN} \leq r_i \leq r_{MAX} \quad \forall i$$

$$c_i = \frac{\sum_{m=1}^M \frac{v_{m,i}}{C_{m,i}}}{\max_i(c_i)} \quad \forall i$$

$$\gamma = \frac{\left(\sum_{i=1}^N d_i \right)}{\sum_{i=1}^N c_i (d_i - r_{i,MIN})^2}$$

$$d_i = \rho_{R,NB} d_{NB} + \rho_{L,SB} (1 - \rho_{R,SB}) d_{SB} + \rho_{T,EB} d_{EB} + \frac{q_i(0)}{T} \quad \forall i$$

which can be solved with any constrained nonlinear programming or specialized quadratic programming method, not detailed here. Note that the inclusion of the overcapacity variables z_p , $i=1 \dots N$ ensures at least the feasible solution

$$r_i = r_{i,MIN}, \quad z_i = d_i T - r_{i,MIN} T - Q_i \quad i = 1 \dots N \quad \text{Eqn. 5- 19}$$

for reasonable (realistic) values of the freeway capacities CAP_j . In the evaluation results of this chapter and Chapter 9, problem (5-18) is solved using the QP barrier algorithm of the CPLEX math programming optimization software package [CPLEX, 1997]. For more details of the iteration details of barrier optimization algorithms, see [Bazaraa et al., 1993].

Operation under severe congestion

If any

$$CAP_j \leq \sum_{i=0}^j A_{i,j} r_{i,MIN} \quad \text{Eqn. 5- 20}$$

then there is a *severe* limitation of capacity (an *incident*) in that section and even (5-18) with the inclusion of the z_i variables will be infeasible. In this case, we can prescribe a heuristic solution such that $r_i = r_{i,MIN}$ for all ramps upstream of the congestion with condition (5-20) and $r_i = r_{i,MAX}$ downstream of the severely congested section.

In addition, the higher-layer processor(s) of the ramp metering control system should send information to the surface-street controllers regarding the incident location (severe limitation of capacity) and the prescribed emergency settings of $r_{i,MIN}$. In the presence of ATIS, such information could also be provided to travelers to increase the diversion effect away from the congested segment. As the congestion clears, the anomaly detection module will detect the favorable change to the state variables x_j and re-run the area-wide optimization for a new, feasible solution to the area-wide coordination problem (5-18).

Integration with predictive-cooperative real-time rate regulation layer

After solving (5-18) to optimality, we obtain the *nominal* ramp metering rates \bar{r}_i , volumes \bar{V}_j , and queue lengths at the end of the time horizon $q_{i,N}(K)$. From this, we derive the steady-state *set-point* densities $\bar{\rho}_j$, speeds \bar{v}_j from the characteristic volume-density and speed-density curves of freeway flow. These set-points, the dual variables λ_k and slack values ε_k of constraints (5-3), (5-4), and (5-6), are provided to the predictive-cooperative real-time rate regulation layer to derive the real-time flow measurements. More detail of how these issues are addressed is provided in Chapter 6.

Preliminary evaluation of the area-wide coordination problem on a small example

In this section, we describe an example problem and the solution results using the quadratic problem formulation (5-18) versus:

- (1) a linear programming formulation without queue storage considerations
- (2) the no-control case, and
- (3) a policy to set ramp metering rates at 600 veh/hr, regardless of demand.

Consider a long (30 km) eastbound two-lane freeway with 5 controllable ramps at 5 km spacing, as shown in Figure 5-2. At the initial freeway entrance, there is considerable external demand. In this example, we assume each lane of the freeway can carry 2000 vehicles per hour.

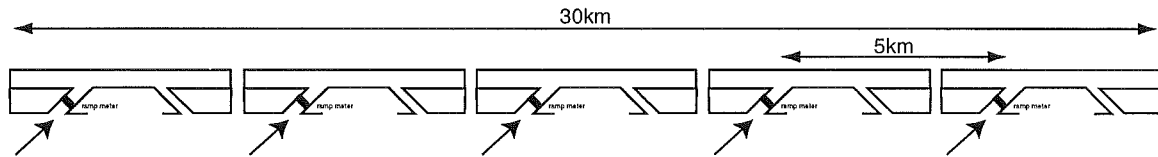


Figure 5- 2. Example problem

This example is similar to the example used in the work by Papageorgiou in the development of a large-scale hierarchical freeway control system [Papageorgiou, 1983]. The parameters and essential data for this example problem are presented in Tables 5-1 and 5-2. Recall from earlier discussion in this chapter that the parameters $A_{i,j}$ correspond to the proportion of vehicle flow from on-ramp i continuing through freeway section j .

from\to	begin	1	2	3	4	5
begin	1	0.95	0.9	0.85	0.8	0.75
1		1	0.95	0.9	0.85	0.8
2			1	0.95	0.9	0.85
3				1	0.95	0.9
4					1	0.95
5						1

Table 5- 1. Route-proportional matrix of example problem

MILOS requires information about the interchange flows that comprise the ramp demand, and needs performance information to derive the objective function cost coefficients. The data in Table 5-2 includes the volumes and turning probabilities at each interchange, that make up the demands to each ramp. Coupling these data with the green-time percentages (GT%), the V/C ratios of each interchange are computed and the scaled cost coefficients c_i are computed from the V/C ratios. These data have been selected to produce values for the demands at each ramp similar to the example problem of Papageorgiou [1983] and, at the same time, to produce differentiation in the congestion level at each interchange. As indicated by the V/C ratio at each location, interchanges 3 and 4 are more congested than 1, 2, and 5. As a result, the QP coordination algorithm should store less vehicles on ramps 3 and 4 (relative to demand at that ramp and the storage capability of the ramp) than at ramps 1, 2, and 5.

The values chosen in Table 5-2 do not represent real locations, but they are intended to be reasonable approximations of real behavior. The interchange parameters have also been selected to illustrate the ability of the quadratic programming formulation to react to differences in the congestion level at adjacent ramps. The quadratic programming problem (5-18) was solved using the data in Tables 5-1 and 5-2 with several values of β , the trade-off parameter between the two sub-objectives. The solutions, as listed in Table 5-6, indicate that the LP method restricts on-ramp volume at ramps 3, 4, and 5 *only*, where the freeway capacity is exceeded. The QP methods, since they take into account the queue storage restriction, restrict on-ramp volume at *all* ramps. The last line of Table 5-3 indicates when the freeway would become over-capacity if *no* ramp metering was implemented. Thus, in this example problem, if no ramp metering was implemented a backward-traveling congestion wave would be started in the section containing ramp 3.

On-Ramp data	begin	1	2	3	4	5
Northbd vol (veh/hr)		2200	1600	1800	1000	1100
Southbd vol (veh/hr)		2000	2000	1900	1000	900
Eastbd vol (veh/hr)		200	260	282	295	307
$P_{nb,r}$		0.25	0.35	0.15	0.25	0.25
$P_{sb,l}$		0.1	0.05	0.07	0.1	0.1
$P_{eb,t}$		0.05	0.07	0.05	0.15	0.05
Ramp demand d_i (veh/hr)	3000	684	610	355	355	342
Initial queue	-	0	0	0	0	0
Time horizon (hrs)		0.5	0.5	0.5	0.5	0.5
Ramp storage (veh)	-	40	30	40	40	50
GT%, EB	-	0.18	0.25	0.15	0.2	0.23
GT%, NB/SB	-	0.55	0.6	0.7	0.55	0.6
GT%, SB left	-	0.27	0.15	0.15	0.25	0.17
v/C ratio		1.67	1.55	1.74	1.24	1.15
scaled cost c_i		0.95	0.89	1	0.71	0.66

Table 5- 2. Ramp interchange data

	ramp 1	ramp 2	ramp 3	ramp 4	ramp 5
Demand	684	610	355	355	342
LP	684	610	320	275	275
QP, $\beta=1, \beta_2=100\beta$	642	555	295	275	253
QP, $\beta=100, \beta_2=100\beta$	644	550	275	275	271
QP, $\beta=0.01, \beta_2=100\beta$	669	550	275	275	252
6-second cycle	600	600	355	355	342
freeway over capacity?	NO	NO	YES	YES	YES

Table 5- 3. Comparison of metering rate coordination methods

Influence of β

In this example, the differences between the QP solutions for radically different values of β are negligible because of the queue-growth constraints (5-6). In addition, setting $\beta_2 = 10000\beta$ in each case prescribes a solution that allows no queues to spillback regardless of their congestion level. It is possible with a lower setting for β_2 (relative to the value of β) that spillback could occur at a ramp with a relatively uncongested interchange. As such, care must be taken in setting the value for β_2 when queue balancing is de-emphasized.

β has the most influence on the resulting rate allocation when metering is required to avoid freeway congestion near interchanges where the congestion is also considerable. In such a situation, an LP approach that does *not* consider queue restrictions or interchange congestion may apply restrictive metering at locations where the most adverse impact on the surface-streets would result. The QP approach of (5-18) would enact metering upstream at interchanges where the congestion was (possibly) lower (i.e. $c_i < c_j \mid i < j$) and resulting ramp queues would have less effect on the surface-street congestion.

To further illustrate that the QP solution balances queues according to interchange congestion, consider Table 5-4. Table 5-4 compares the rates for each ramp solved by the QP method for $\beta=1$, $\beta_2=100\beta$. The queue storage size Q_i , demand rate d_i , A_{ij} , and cost coefficient c_i of each interchange all interact to produce the metering rates r_i . Comparing the rate at ramp 5 from the QP and LP in Table 5-3 provides some evidence that more vehicles are held at uncongested ramps with the QP approach. Table 5-4 confirms this, since ramp 5 also has the largest ramp storage capacity Q_5 .

	ramp 1	ramp 2	Ramp 3	ramp 4	ramp 5
demand (veh/hr)	684	610	355	355	342
Metering rate (veh/hr)	642	555	295	275	253
Vehicle storage size	40	30	40	40	50
Cost coefficient	0.95	0.89	1	0.71	0.66
Percentage holdback	6.1	9.0	16.9	22.5	26.0
Maximum queue, 30-min horizon	21	28	30	40	45

Table 5- 4. Rate comparison for $\beta=1$

Comparison of area-wide metering rate settings in macrosimulation

We now compare the area-wide coordination optimization algorithm based on quadratic programming (5-18) with the no-control and LP-control cases in the macrosimulation environment of Chapter 4. For the example problem of Figure 5-2, we use parameters of the macrosimulation similar to those used by Papageorgiou [1984] with the addition of the cut-off level ρ_c which phases in the flow limitation in (4-11) and (4-12). These parameters are listed in Table 5-5.

0.013	21.6	10	2	4	0.004 hr	0.8	100	123	80
τ	v	κ	l	m	T	α	ρ_{Max}	v_{Max}	ρ_c

Table 5- 5. Parameters for example problem

Each of the area-wide rate coordination solutions listed in Table 5-3 were tested using the macrosimulation environment. Each simulation was run with the initial conditions listed in Table 5-6. The initial speeds $v_j(0)$ were computed from the theoretical speed-density relationship $v_j = v_e(\rho_j)$ of (4-5). Each simulation was executed for one hour of simulated time using the constant route proportional matrix A given in Table 5-1 and the constant input demands d_i given in Table 5-2.

state variable / section	section 1	section 2	section 3	section 4	section 5
density (veh/km)	51	59	65	65	65
speed (km/hr)	69.7	68	70.5	70.	70.5
speed limit (km/hr)	123	111	98.4	98.4	98.4

Table 5- 6. Initial conditions for simulation run

Figure 5-3 depicts the time-space plots of density through the corridor over the duration of the simulation for the no-control, LP-control, and QP-control case ($\beta=1$, $\beta_2=1$), respectively. Figure 5-4 illustrates the differences in the queue growth using the LP method and the QP method. Queue growth for the no-control case is not shown in Figure 5-4 but it should be noted that a small queue of 6 vehicles develops and dissipates at the second on-ramp when the backward-traveling congestion wave passes the ramp around the 40-minute mark of the simulation. This queue buildup is due to the inclusion of the soft-limiter in (4-12) and would *not* have been represented in previous versions (prior to this research) of the macroscopic simulator.

Table 5-7 lists the system performance measures of total travel time, queue delay time, and throughput for each of the methods over the duration of the simulation. Throughput is computed by totaling the vehicles leaving the freeway over the duration of the simulation and dividing by the total number of vehicles entering the simulation. Because the simulation ends with vehicles still in the system, the throughput will not be close to unity. Notice, however, that the QP method provides approximately 2% more throughput than the LP method and 4% more throughput than no control during the transient. Although the LP method results in lower freeway travel time, the QP method has lower total queue time, even though more ramps are metered. As indicated in Figure 5-4, the LP method spills-back two ramps and builds very small queues at two ramps. The corresponding QP

solution does not result in queue spillback at any ramp, and results in queues that are more balanced throughout the corridor.

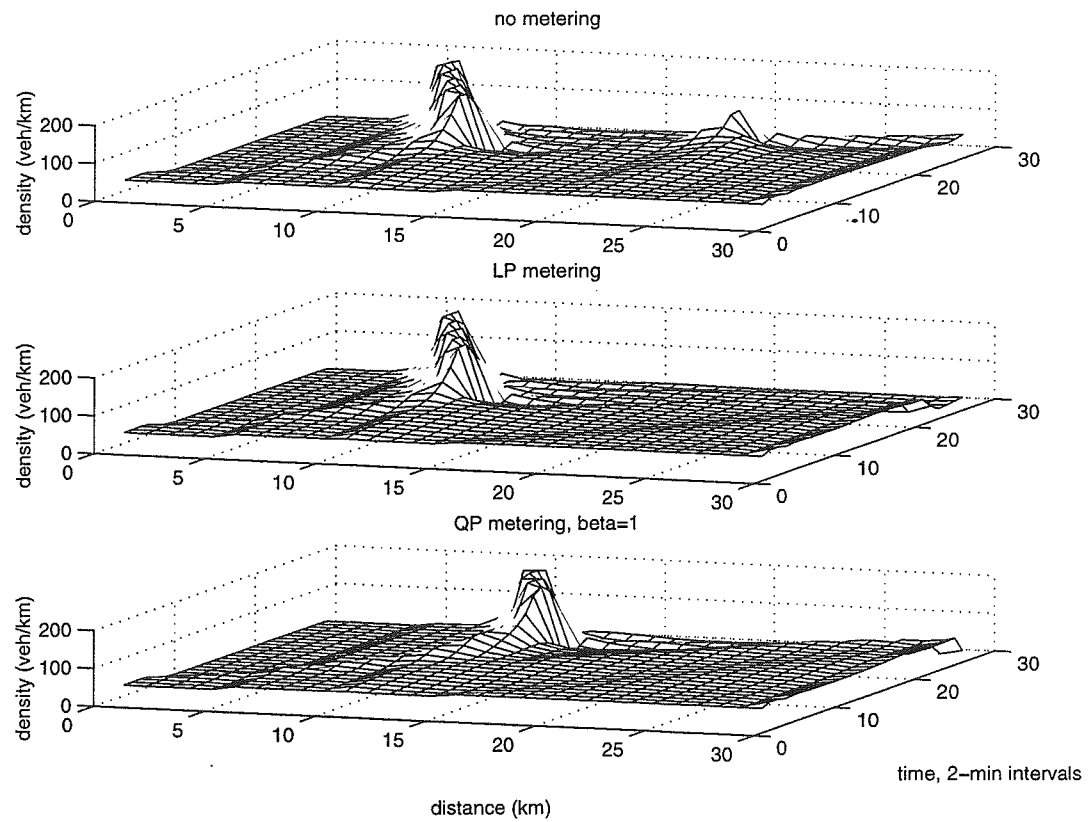


Figure 5- 3. Density evolution comparison

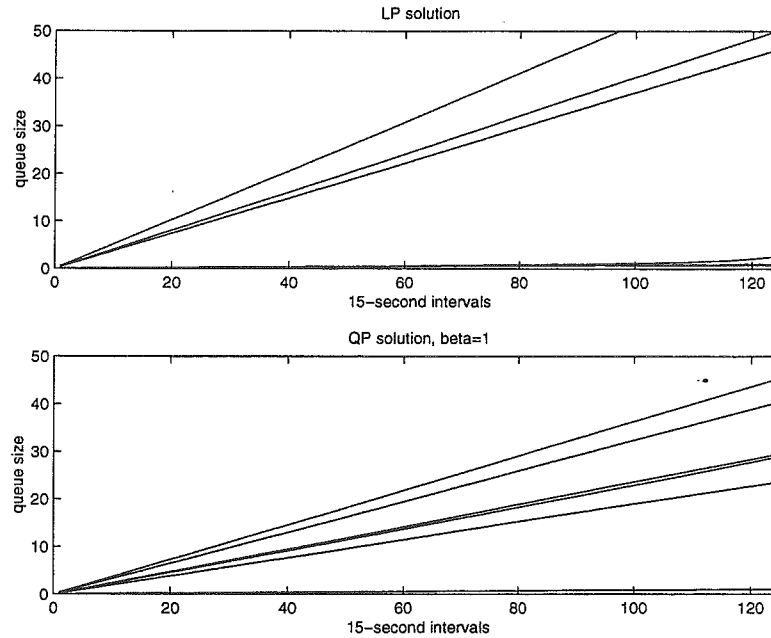


Figure 5- 4. Queue growth comparison

Method / Measure	Fwy travel time (veh-hr)	Queue time (veh-hr)	served load / offered load
No control	989	0.57	3048/5366 = 0.568
LP control	964	144	2976/5046 = 0.589
QP control	941	131	3077/5037 = 0.611

Table 5- 7. Preliminary method comparisons

Simulation test with extended queue dissipation

The previous comparison may be somewhat biased towards indicating the LP and QP methods are superior because the simulation ends with a significant number of vehicles still stored in ramp queues. Thus, another deterministic simulation was executed for a total time of two hours. The first 30 minutes was run with a set of low input volumes d_i , not requiring ramp metering, then 30 minutes of the high volumes from Table 5-8, and then an additional hour of simulation, at lower volumes as indicated in Table 5-8.

Table 5-9 indicates the initial conditions used in this experiment. This simulation would allow the queues that were built at the ramps by using the LP and QP metering methods to dissipate and result in a more *just* comparison of the LP, QP and no control situations. The hypothesis being that the discharged queues could possibly create a secondary congestion when released simultaneously (by any method), where in the no-control case, this

secondary congestion could not be built up because no such queues were created in minutes 31-60 of the experiment.

time period / input stream	external	ramp 1	ramp 2	Ramp 3	ramp 4	ramp 5
0 - 30 minutes	2500	484	410	255	255	242
31 - 60 minutes	3000	684	610	355	355	342
61 - 120 minutes	2500	484	410	255	255	242

Table 5- 8. Demand volumes for evaluation example 2

state variable / section	section 1	section 2	section 3	section 4	section 5
density (veh/km)	18	21	26	26	26
speed (km/hr)	101	101	95	95	95
speed limit (km/hr)	123	111	98.4	98.4	98.4

Table 5- 9. Initial conditions for evaluation example 2

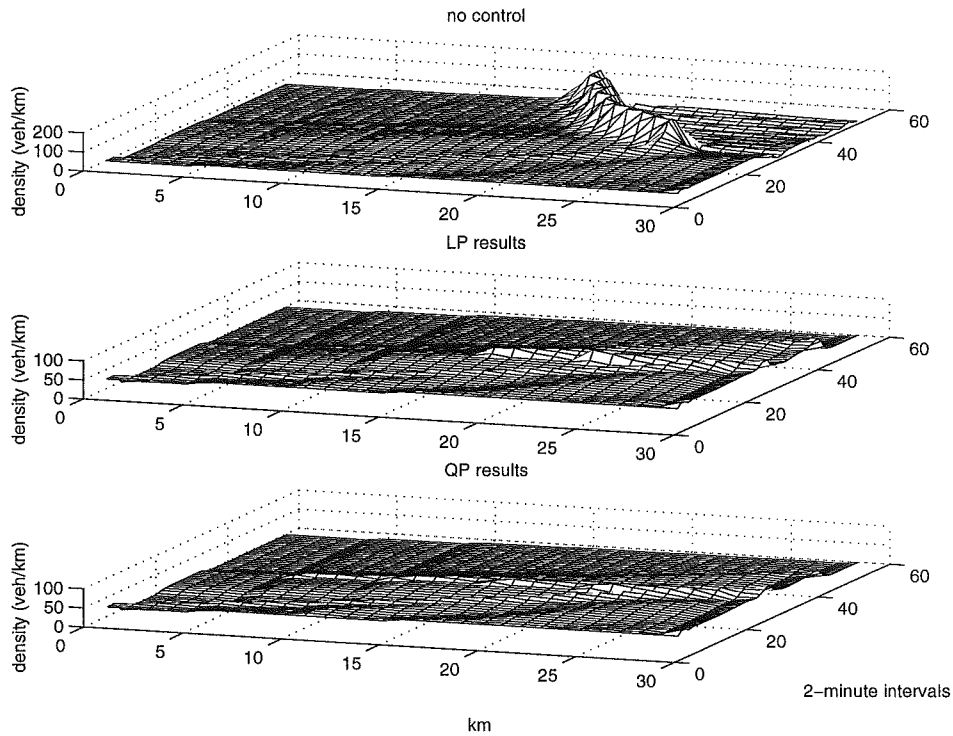


Figure 5- 5. Density evolution comparisons for evaluation example 2

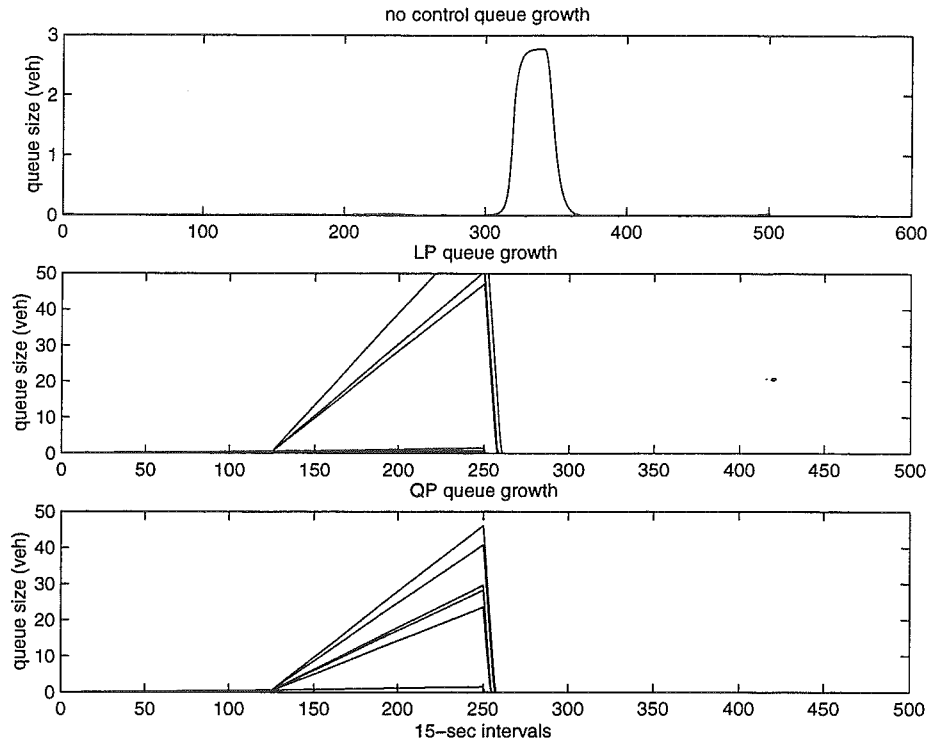


Figure 5- 6. Queue growth comparisons, evaluation example 2

Method / Measure	Fwy travel time	Queue time	offered load/served load
No control	2357	0.35	$4847/9000 = 0.539$
LP control	2262	44.9	$5559/8917 = 0.623$
QP control	2249	45.2	$5559/8917 = 0.623$

Table 5- 10. Performance comparisons, evaluation example 2

In this evaluation, the LP and QP methods have identical throughput, both 9% higher than the no-control case as listed in Table 5-10. The QP method results in slightly less freeway travel time than the LP method. With constant, deterministic arrival rates, this difference is considered statistically insignificant. As indicated in Figure 5-5, the backward-traveling congestion wave of the no-control situation is eliminated using both the LP and QP methods. Notice also in Figure 5-6, as in the previous example problem, that a queue of two vehicles is created and dissipates at the second on-ramp as the backward-traveling congestion wave passes the ramp area in the “no-control” case. This indicates the effect of adding the soft-limiter to queue-length computation in (4-12).

Summary

Overall, this simulation test showed that over a single time-horizon (i.e. without a rolling-horizon implementation) *queue-growth management* is the main benefit provided by the QP area-wide coordination algorithm over the LP approach. In an example problem, the performance of the freeway system under the two methods was virtually identical but two queue limitations were exceeded using when LP was used. The resulting benefits of using the QP approach *off of* the freeway must be inferred since we do not simulate the operation of the interchanges. We conjecture, however, that since the LP method spills-back more queues than the QP method, it is unlikely that the performance of the LP solution at the interchanges is superior to that of the QP result. In addition, as the congestion levels of the interchanges vary, the QP area-wide coordination optimization method reacts to these changes and should show even more powerful benefits. More extensive simulation testing, using stochastic input flows and the full rolling-horizon, two-layer optimization structure of MILOS is detailed in Chapter 9. Next, in Chapter 6, a predictive-cooperative real-time control layer will be described that realizes additional performance benefits (i.e. travel-time savings and queue dissipation) by modifying the nominal ramp metering rate according to the observed stochastic disturbances in the freeway and ramp demand flows.

Chapter 6: Predictive-cooperative real-time rate regulation algorithm

Introduction

The concept of a multi-layer controller structure and use of a regulation algorithm to further refine *nominal* ramp metering rates defined by a higher-level processor is not a new one [Payne and Isaksen, 1973; Papageorgiou, 1983; Payne et al., 1985]. The approach taken in this research is distinguished by several characteristics

- (1) in contrast to the *reactive* linear-quadratic-regulator (LQR) approach based on automatic control theory, we develop a predictive, anticipatory approach that is proactive to *future* traffic states,
- (2) queue management is *explicitly* considered in the real-time optimization formulation,
- (3) a method is presented to *continue* real-time optimization even in the case where one or more segments of the freeway are congested, and
- (4) information about the *dual* of the upper-layer area-wide coordination problem is used to guide real-time subproblem optimization.

Virtually no traffic-reactive ramp metering methods have been developed which *explicitly* consider queue management in the optimization procedure. There have been some preliminary attempts to manage queues at ramps, such as by synchronizing ramp metering rates with surface-street signal phases [Han and Reiss, 1995]. Another report lists "strategies" and "tactics" for integrated freeway/surface-street traffic management [Pooran et al., 1992].

Queue management has been historically addressed by traffic engineers using a *spillback* detector just before the surface-street interchange which create oscillations between high $r_{i,MAX}$ and low rates r_i with their switching behavior [ADOT, 1997]. Recently, a queue-control method has been presented that reduces the oscillating behavior by estimating occupancy on shorter time-scales using a first-order filter [Gordon, 1996]. This algorithm, however, does *not* derive/suggest the "higher" rate to change to when the queue length is to be reduced.

Adding queue management to state feedback control methods

To date, feedback control methods offer the most theoretically-founded solution to the problem of reacting to the inevitable stochastic irregularities in freeway traffic flow. These algorithms compensate for the stochastic disturbances and drive the measured state ρ_j towards a desired "nominal" state $\rho_{j,N}$. This nominal state is provided by an external system, algorithm, or decision-maker [Papageorgiou, 1983]. However, there are several drawbacks of applying linear feedback control methods.

First, when the disturbance is severe, linear feedback control methods may not drive the system back to the nominal state $\rho_j(k) \rightarrow \rho_{j,N}$ because the linear approximation to the system dynamics is no longer valid [Papageorgiou, 1983; Payne et al., 1985]. This can be addressed by turning "off" the local feedback controller in the presence of congestion and re-solving the upper-layer area-wide coordination optimization problem with a reduced-flow restriction in the congested section(s). When the system has returned to normal operation, the local feedback controller can then be turned back "on". We will show later that it is possible to continue locally traffic-responsive rate regulation even in the presence of congestion by using the predictive-cooperative real-time rate regulation algorithm developed in this chapter. The second main drawback of linear feedback control methods is the difficulty of such methods to consider queue management *explicitly* in the optimization procedure. These methods do not appropriately model the costs of changes to the queue and the matrix ricatti equation cannot be solved since in the system model the state-variable coefficient matrix A_c is singular.

Central concept of PC-RT rate regulation algorithm

The method developed in this chapter for queue management is termed the predictive-cooperative real-time rate regulation algorithm, hereafter referred to as PC-RT. The PC-RT rate regulation algorithm addresses the need to integrate the control effort of freeway control system with the concerns of the surface-street control system by:

- (a) responding to statistically-significant short-term fluctuations in the stochastic demand flows to the ramp system (i.e. *both* the upstream freeway flow and the ramp demand),

- (b) satisfying the coordination requirements of the area-wide problem solved at the upper layer,
- (c) scheduling rates that are *pro-active* to possible future ramp and freeway demands, and
- (d) reducing the possibility of and managing queue spillback at ramp entrances.

Point (a) reflects our assumption that, at the area-wide coordination layer, ramp demands and upstream freeway flows are considered constant over a short (i.e. 15-minute) time horizon. However, during that time horizon, it is well known that demand is *not* constant due to stochastic fluctuations. As such, any constant metering rate $r_i(k) = r_{i,N} \forall k$ will neglect these fluctuations and create queues when not necessary and/or release vehicles when disadvantageous to freeway conditions. It is important, however, at the real-time layer to distinguish between *negligible* statistical variation in the traffic stream and *significant* deviations in the flow. An algorithm for identifying these deviations based on statistical process control is discussed in Chapter 7.

Point (b) indicates that the PC-RT rate regulation algorithm should continue to satisfy the *area-wide* metering objectives by applying rates that do not deviate significantly from the nominal rates recommended by the upper-layer coordination optimization problem that is, $r_i(k) \in [r_{i,N} \pm \Delta r_{i,N}]$. One of the issues to be evaluated in this research is then, obviously, can a ramp metering control system be both coordinated on an area-wide basis, yet remain locally traffic-responsive?

To provide a preliminary response to this issue, consider the analogy of the *cycle, split, offset* concept from surface street control. In this concept, network coordination is maintained by using pre-timed signal controllers that operate on the same background cycle length and setting the splits and offsets to provide progression opportunities along arterials [McShane and Roess, 1990]. With the same cycle, split, offset paradigm, *semi-actuated* signal controllers can be used in the coordinated surface-street control system to provide some reactivity to local traffic conditions at an intersection while at the same time maintaining the progression opportunities of the coordinated system.

Finally, points (c) and (d) describe the *predictive* and *cooperative* components of the PC-RT rate regulation algorithm, respectively. The PC-RT rate regulation algorithm is designated as *predictive* to indicate that the control algorithm evaluates several *scenarios* for a short (i.e. five to seven minute) time horizon into the future of what *could* happen to the upstream freeway flow $V_o(k+1)$, $V_o(k+2)$, ..., $V_o(k+N)$ and ramp demands $d_i(k+1)$, $d_i(k+2)$, ..., $d_i(k+N)$. The PC-RT rate regulation algorithm is designated as *cooperative* to indicate the focus of the algorithm on *cooperating* with the surface-street signal controller to manage queue spillback by adapting the metering rate to react to the next few minutes of predicted flows from the interchange. Hence, the *cooperative* queue management method proposed in the PC-RT rate regulation algorithm does *not* default to the oscillating behavior resulting from using an occupancy threshold at a static “queue detector” to raise the metering rate to some pre-set “high” rate when the queue is sufficiently long.

Anticipated effects of PC-RT rate regulation algorithm

At the very least, the intended effect of implementing the PC-RT rate regulation algorithm with the solution from the QP area-wide coordination optimization problem should make the overall system (i.e. freeway and ramp queue) performance that is:

- (a) *equivalent to or better than* the area-wide QP solution alone,
- (b) *equivalent to or better than* feedback control methods [Papageorgiou, 1984, 1991; Payne et al., 1985] which do *not* consider queue management in their optimization procedures, and
- (c) *equivalent to or better than* volume-occupancy traffic-responsive metering algorithms.

Thus, the intent of this research is to show that queue management can be considered in deriving metering rates that result in similar, if not better, freeway and surface-street performance than methods that do not consider queue management explicitly.

Basic function of the PC-RT rate regulation algorithm

The basic function of the PC-RT rate regulation algorithm is to exploit, at any time k , the excess local capacity $\rho_j(k) < \rho_{j,N}$ and $q_i(k) < q_{i,N}(k)$ in the freeway/ramp system by reacting in the following ways to the *fundamental* combinations of predicted ramp demand and predicted upstream freeway flow:

- (1) *increase* the metering rate when the freeway density is *lower* than the nominal density and the ramp demand is *higher* than nominal,
- (2) *decrease* the rate when the ramp demand is *lower* than nominal and freeway density is *higher* than nominal
- (3) *increase* the rate when ramp demand is *lower* than nominal and freeway density is *lower* than nominal
- (4) *increase* or *decrease* the metering rate according to a *trade-off* solution when ramp demand is *higher* than nominal and freeway density is *higher* than nominal.

How much to decrease or increase the rate $r_i(k)$ from the nominal setting $r_{i,N}$ is specified by formulation of a linear programming optimization problem (LP). This LP is formulated with a linearized description of the macroscopic freeway flow equations (from Chapter 4) about the nominal *equilibrium state* $(\rho_{j,N}, v_{j,N}, r_{i,N})$ and a linear description of queue growth about the nominal queue-growth *trajectory* $q_{i,N}(k)$. The cost function of this LP optimization problem is a weighted sum of travel-time savings in each section of the freeway and on the ramp approaches. The weights of each state-variable are derived from the dual multipliers λ_k and constraint slack ε_k values of the solution to the upper-layer area-wide QP optimization problem. In this manner, a trade-off between travel-time savings on the ramp and on the freeway is based on the *current* interchange conditions (i.e. how important it is to manage spillback at *this* ramp) and the conditions in critical freeway sections.

The PC-RT rate regulation algorithm can be described as a three-step process:

- (1) Given that a *significant* deviation from the upstream freeway or ramp demand nominal flow is detected, *predict* several possible subsequent flows to the ramp and the upstream freeway segment,
- (2) Given these predicted possible future *scenarios*, solve an LP optimization problem for each predicted scenario that reduces queuing time on the ramp and/or reduces the possibility for congestion on the freeway over the next few minutes, and
- (3) In the next optimization interval, collect the "actual" upstream freeway flow and ramp demand, compare the actual flow to the predicted scenarios, and apply the appropriate metering rate for the scenario that best matches the actual flow.

A *rolling-horizon* framework is used in the three-step process listed above. Thus, the PC-RT optimization problems are solved for a 5 to 7 minute predictive time-horizon, but the metering rate is only applied for the first 1-2 minutes of the time horizon before the problem is possibly re-evaluated due to the stochastic fluctuations.

Reasonable and important assumptions

Some reasonable and important assumptions are needed to successfully formulate and apply meaningful optimization results at the predictive-cooperative, real-time rate regulation layer:

- (1) The traffic flows to ramps and on the freeway have a stochastic component that can be *identified* and *separated* from structural changes in the underlying process. For the case of ramp demands, it is assumed that this stochastic effect can be separated from the flow-rate changes induced by the traffic signal at the interchange.
- (2) At least one-minute upstream measurements of flows to the ramp meter are available from the surface street detector system and are reliable enough for use in the PC-RT rate regulation and area-wide coordination optimization problems.
- (3) The linear approximation of the nonlinear macroscopic flow equations (Chapter 4) can be used to accurately predict freeway flow dynamics over 1-7

minute time scales. Correspondingly, route-proportional flow rates and other parameters of the linearized model are assumed to remain constant over the real-time optimization time horizon.

- (4) The flows and metering rates solved for at the area-wide coordination layer of approximation are realizable as an equilibrium-state, in the absence of stochastic variation.
- (5) The flows resulting from solution of the area-wide coordination problem are *desirable* settings from a system-wide perspective. That is to say that the solution of the area-wide problem provides benefits that *individual* solution of the local control problems, with independently-derived nominal settings (e.g. critical flows that do not vary by time-of-day), could not provide. Thus, the goal of keeping the real-time rates "close to" the area-wide nominal settings is a sound objective.

Linearization about an equilibrium state

Having an equilibrium solution and a set of desired nominal rates that address network-wide concerns from the area-wide coordination problem, we can begin to simplify the difficult optimal control problem to one that can be solved in real-time. First, as developed in previous work [Papageorgiou, 1983; Payne et al., 1985], consider the differential equation model for the freeway density dynamics $\dot{\rho} = f(\rho, v, r, d, t)$. ρ is considered the *primary* dynamic variable in freeway flow modeling since, as shown by empirical data, speeds are approximately constant until the density approaches the critical value ρ_c . As such, the dynamic equation for the speed v is replaced by the equilibrium speed-density function $v_e(\rho)$ in the freeway flow dynamics for the development of the PC-RT optimization subproblems.

Hence, we approximate the nonlinear function $f(\rho, v, r, d, t)$ by a linearization (first-order Taylor-series expansion) about the equilibrium point (ρ_N, v_N, r_N)

$$\hat{\rho}_j = f_j(\rho_N, r_N) + \sum_{i=1}^J \left. \frac{\partial f_i(\rho, r)}{\partial \rho_i} \right|_{\rho_N, r_N} (\rho_j - \rho_{j,N}) + \sum_{j=1}^J \left. \frac{\partial f_j(\rho, r)}{\partial r_j} \right|_{\rho_N, r_N} (r_j - r_{j,N}) \quad \text{Eqn. 6- 1}$$

where J indicates the number of segments in the freeway model. Each segment of the freeway may not contain an off-ramp, hence many $r_j = 0$. We include those terms for completeness only. A linear system description results for the deviation from the nominal state

$$\Delta \dot{\rho}_j = \mathbf{A}^* \Delta \rho_j + \mathbf{B}^* \Delta r_j \quad \text{Eqn. 6- 2}$$

such that

$$\begin{aligned} \Delta \rho_j &= \rho_j - \rho_N \\ \Delta r_j &= r_j - r_N \\ \Delta \dot{\rho}_j &= \dot{\rho}_j - \dot{\rho}_N \\ \dot{\rho}_N &= f(\rho_N, r_N) \end{aligned} \quad \text{Eqn. 6- 3}$$

and

$$\mathbf{A}^* = \begin{bmatrix} \left. \frac{\partial f_1(\rho, r)}{\partial \rho_1} \right|_{\rho_N, r_N} & \dots & \left. \frac{\partial f_1(\rho, r)}{\partial \rho_J} \right|_{\rho_N, r_N} \\ \vdots & \ddots & \vdots \\ \left. \frac{\partial f_J(\rho, r)}{\partial \rho_1} \right|_{\rho_N, r_N} & \dots & \left. \frac{\partial f_J(\rho, r)}{\partial \rho_J} \right|_{\rho_N, r_N} \end{bmatrix}, \quad \mathbf{B}^* = \begin{bmatrix} \left. \frac{\partial f_1(\rho, r)}{\partial r_1} \right|_{\rho_N, r_N} & \dots & \left. \frac{\partial f_1(\rho, r)}{\partial r_J} \right|_{\rho_N, r_N} \\ \vdots & \ddots & \vdots \\ \left. \frac{\partial f_J(\rho, r)}{\partial r_1} \right|_{\rho_N, r_N} & \dots & \left. \frac{\partial f_J(\rho, r)}{\partial r_J} \right|_{\rho_N, r_N} \end{bmatrix} \quad \text{Eqn. 6- 4.}$$

with $f_j(\rho, v, r)$ defined as

$$f_j(\rho, v, r) = V_{j,IN} - V_{j,OUT} + r_j - s_j \quad \text{Eqn. 6- 5}$$

and

$$\begin{aligned} V_{j,IN} &= \alpha v_{j-1} \rho_{j-1} + (1 - \alpha) v_j \rho_j \\ V_{j,OUT} &= \alpha v_j \rho_j + (1 - \alpha) v_{j+1} \rho_{j+1} \\ s_j &= \theta_j v_j \rho_j \end{aligned} \quad \text{Eqn. 6- 6}$$

as developed in Chapter 4. Collecting like terms in (6-5) after substituting the relationships of (6-6) we get

$$f_j(\rho, v, r) = \alpha v_{j-1} \rho_{j-1} + \{(1 - 2\alpha - \theta)\} v_j \rho_j - (1 - \alpha) v_{j+1} \rho_{j+1} + r_j \quad \text{Eqn. 6-7}$$

This formulation uses the *turning-percentage* representation of the off-ramp rate, indicating that, as the density and speed in the section change, so does the number of vehicles exiting the section. An alternative to this proportional off-ramp rate model is to use the (fixed) off-ramp rate specified by the nominal solution of the upper-layer area-wide coordination problem.

Elimination of the dynamic speed equation

The complex nonlinear dynamic equation (4-3) for v_j is replaced by solving for the equilibrium speed $v_{j,N}$ from the nominal density $\rho_{j,N}$ from (4-5)

$$v_{j,N} = v_e(\rho_{j,N}) = v_f \left(1 - \left(\frac{\rho_{j,N}}{\rho_{j,MAX}} \right)^t \right)^m \quad \text{Eqn. 6-8}$$

and thus the matrices A^* and B^* can hence be written as

$$A^* = \begin{bmatrix} a_{1,1} & a_{2,1} & 0 & \dots & 0 \\ a_{1,2} & a_{2,2} & a_{3,2} & & \dots \\ 0 & & \dots & & 0 \\ \dots & & a_{M-2,M-1} & a_{M-1,M-1} & a_{M,M-1} \\ 0 & \dots & 0 & a_{M-1,M} & a_{M,M} \end{bmatrix} \quad \text{Eqn. 6-9}$$

$$B^* = \begin{bmatrix} 1 & & & & \\ & 1 & & & \\ & & \dots & & \\ & & & 1 & \\ & & & & 1 \end{bmatrix} \quad \text{Eqn. 6-10}$$

where elements of A^* are defined as

$$\begin{aligned}
a_{j,j} &= \left. \frac{\partial f_j(\rho, r)}{\partial \rho_j} \right|_{\rho_N, r_N} = \{(1 - 2\alpha - \theta)\} [v_e(\rho_{j,N}) + v_e'(\rho_{j,N})\rho_{j,N}] \\
a_{j,j-1} &= \left. \frac{\partial f_j(\rho, r)}{\partial \rho_{j-1}} \right|_{\rho_N, r_N} = \alpha [v_e(\rho_{j-1,N}) + v_e'(\rho_{j-1,N})\rho_{j-1,N}] \\
a_{j,j+1} &= \left. \frac{\partial f_j(\rho, r)}{\partial \rho_{j+1}} \right|_{\rho_N, r_N} = -(1 - \alpha) [v_e(\rho_{j+1,N}) + v_e'(\rho_{j+1,N})\rho_{j+1,N}]
\end{aligned} \tag{Eqn. 6- 11}$$

where $v_e'(\rho_{j,N})$ is the derivative of $v_e(\rho_{j,N})$ at the nominal density $\rho_{j,N}$. Recall that α is calibrated to the freeway location and θ is the percentage of volume exiting at the off-ramp. B^* can be specified as a diagonal unity matrix because $r_j=0$ for sections that do not have an on-ramp.

The replacement of the nonlinear model with its linearization results in significant reduction of the model order (by elimination of the v state-variables for section speeds) without significant loss of descriptive accuracy around the nominal state point (ρ_N, v_N, r_N) as developed in previous research [Papageorgiou, 1983; Payne et al., 1985].

By using the (simple) definition of the derivative

$$\frac{\rho(t + \Delta t) - \rho(t)}{\Delta t} \tag{Eqn. 6- 12}$$

the continuous nonlinear dynamic system around the equilibrium state is converted to the difference equation

$$\Delta \rho(t + \Delta t) = \Delta \rho(t) + \Delta t (\mathbf{A}^* \Delta \rho + \mathbf{B}^* \Delta r_{on}) \tag{Eqn. 6- 13}$$

for simulation on a digital computer and formulation of the PC-RT rate regulation problems as an LP.

The range of descriptive accuracy of the linear approximation to the nonlinear model has not been studied empirically, but analytical studies have indicated that the equilibrium point (ρ_N, v_N, r_N) is *not* globally asymptotically stable [Papageorgiou, 1983; Zhang and

Ritchie, 1996]. That is, a significant disturbance can drive the system away from the equilibrium point (ρ_N, v_N, r_N) and into a congested “equilibrium” state. This limitation of the region of attraction of the equilibrium point (ρ_N, v_N, r_N) is used as our first set of constraints for the PC-RT optimization problem. Thus, the state variables $\Delta\rho_j$ must remain within the limits of modeling accuracy and applicability such that

$$\max[\rho_{j,MIN}, -\rho_{j,N}] \leq \Delta\rho_j \leq \min[\rho_{crit} - \rho_{j,N}, \rho_{j,MAX}]. \quad \text{Eqn. 6- 14}$$

The lower bound of constraint 6-14 results from the fact that the density must be non-negative, but the model may not be applicable for densities anywhere near zero (when $\rho_{j,MIN} >> 0$). The upper bound of (6-14) results from the system description being invalid for density values greater than the critical value ρ_{crit} , but may still only be applicable up to some point less than the critical point $\rho_{j,N} + \rho_{j,MAX} \ll \rho_{crit}$. In this research, however, we use bounds determined by engineering judgment. Identification of analytical form for these modeling limitations may be a subject of future research.

Similarly, from a modeling perspective the metering rates Δr_j are restricted to the set

$$r_{j,MIN} - r_{j,N} \leq \Delta r_j \leq \hat{r}_{j,MAX} - r_{j,N} \quad \text{Eqn. 6- 15}$$

such that $r_{j,MIN}$ and $r_{j,MAX}$ may be set to some value higher than (lower than) the absolute lowest (highest) metering rate to maintain the model’s accuracy to real system behavior. The constraints 6-15 are also needed to satisfy the coordination requirements of the upper-layer QP problem. Hence, the real-time rates r_j should remain close to the nominal rate. Currently, the $r_{j,MIN}$ and $r_{j,MAX}$ bounds are set proportional to the nominal metering rate, such that

$$\begin{aligned} r_{j,MIN} &= \omega r_{j,N} \\ r_{j,MAX} &= \frac{r_{j,N}}{\omega} \end{aligned} \quad \text{Eqn. 6- 16}$$

where $\omega \in (0,1)$ is a parameter to be chosen by the system operator or based on the attraction region of the equilibrium point (ρ_N, v_N, r_N) .

Structure of the PC-RT objective function

In the PC-RT rate regulation problem, we are interested in adding the objective of managing queue lengths (and minimizing them when possible) to the standard objective of maintaining smooth freeway flow and minimizing total travel time. We use a *linear* objective of the total "additional" travel time in the system as the overall objective of the PC-RT rate regulation problem, since our variables are now defined as the *deviation* from the nominal state.

This implies that, with the state at the nominal point (ρ_N, v_N, r_N) , a certain total travel time

$$\sum_{k=1}^K \left[\sum_{j=1}^N \Gamma_j \rho_{j,N}(k) + \sum_{i=1}^M q_{i,N}(k) \right], \text{ where } \Gamma_j \text{ is the length of segment } j \text{ and the real ramp}$$

locations are indexed by i , would be realized for the system over the time horizon. However, because of stochastic variation in the ramp demands $d_i(k)$ and upstream freeway demand $V_{o,N}(k)$, the travel time will be higher or lower than the expected travel time according to this stochasticity. Our objective in the PC-RT rate regulation problem is to minimize a weighted combination of any *additional* travel time incurred by stochastic variation, and, at the same time, exploit the stochastic fluctuations to our advantage by dissipating queues at appropriate times. Hence, we define the objective of the PC-RT problem as

$$\min \sum_{k=1}^K \left[\sum_{j=1}^J c_j \Gamma_j \Delta \rho_j(k) + \sum_{i=1}^M c_{i,q} \Delta q_i(k) \right] \quad \text{Eqn. 6- 17}$$

where c_j is a weighting cost coefficient of freeway section j , $\Delta \rho_j(k)$ is the deviation from the nominal point $\rho_{j,N}$ of freeway section j at time k , $c_{i,q}$ is a weighting coefficient of queue i , and $\Delta q_i(k)$ is the deviation from the nominal queue length $q_{i,N}(k)$ of the queue at ramp i at time k . Now we denote the total number of ramps in the problem as M . Derivation of the cost coefficients c_j and $c_{i,q}$ will be detailed in a later section.

The first component of the cost function (6-17) reflects the desire to keep the freeway from being forced into congestion by the stochastic variation in the flow rate. Hence, at the expense of throughput (i.e. higher densities), we would like to keep the density on the freeway at or below the nominal point $\rho_{i,N}$ (i.e. $\Delta\rho_i(k) < 0$) specified by the upper-layer area-wide coordination problem. This reduction in density is achieved by holding (even more) vehicles back on the ramp, which is in direct conflict with the concerns of the surface-street control system to keep the ramp queue from spilling back into the adjacent interchange.

Queue growth modeling

The second part of the PC-RT objective function (6-17) reflects the queue-management concerns of the surface-street system where $c_{i,q}$ is the cost coefficient reflecting the *importance* of deviations of the queue length $q_i(k)$ at queue i from the *nominal* queue $q_{i,N}(k)$ length at time k . We now define the concept of a “nominal queue”. Consider that, at the area-wide optimization layer of control, we specify ramp metering rates $r_{i,N}$ that are based on a constant arrival rate d_i at each ramp. Applying these constant rates $r_{i,N}$ with constant demands d_i results in a (known) queue length $q_i(k)$ when the initial queue length $q_i(0)$ is known.

Hence, $q_{i,N}(k)$ could be considered the “acceptable” queue as a limit of how long we would be willing to allow the queue to grow to at each time instant k over the time horizon K . Thus, the deterministic queue growth function $q_{i,N}(k)$ becomes the nominal *trajectory* of the new state variable

$$\Delta q_i(k) = q_i(k) - q_{i,N}(k) \quad \text{Eqn. 6- 18}$$

with the evolution equation

$$\Delta q_i(k+1) = \Delta q_i(k) + \Delta T[\Delta d_i(k) - \Delta r_i(k)] \quad \text{Eqn. 6- 19}$$

such that $\Delta d_i(k)$ is defined as the deviation of the current ramp demand from the constant demand assumed at the area-wide optimization layer

$$\Delta d_i(k) = d_i(k) - d_{i,N} \quad \text{Eqn. 6- 20}$$

and $\Delta r_i(k)$ is defined as the deviation of the current metering rate $r_i(k)$ from the constant rate $r_{i,N}$ specified from the area-wide optimization problem such that

$$\Delta r_i(k) = r_i(k) - r_{i,N}. \quad \text{Eqn. 6- 21}$$

Coupling the definition (6-18) of the new variable $\Delta q_i(k)$ with the objective of minimizing this variable, it should be apparent that the PC-RT objective (6-17) reflects minimizing the *additional* travel time in the system. Hence, the goal of the PC-RT controller is to *avoid* additional congestion in the freeway system while taking advantage of *opportunities* to release vehicles from ramp queues. This situation is illustrated in Figure 6-1.

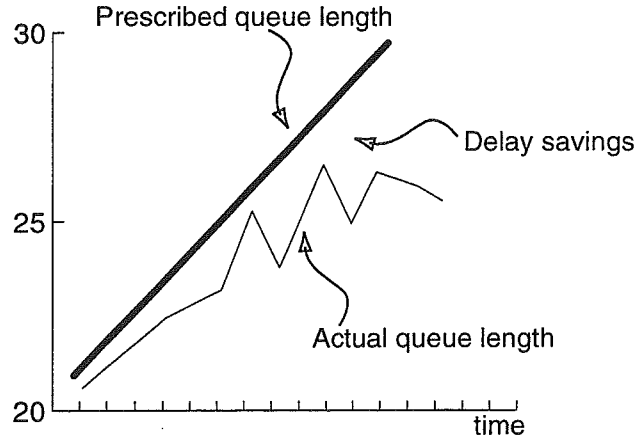


Figure 6- 1. Prescribed maximum queue growth rate

As shown in Figure 6-1, we allow the queue to grow from $q_i(0)$, at most, on average, as fast as the deterministic rate $d_{i,N} - r_{i,N}$ specified in the solution to the area-wide coordination problem. This is illustrated in Figure 6-1 as the thick dark line. Thus, at the end of the short-term time-horizon of the PC-RT optimization, the queue length should be no greater than the length specified by applying the constant rate $r_{i,N}$. This does not mean that faster rates of queue growth are not allowed, but, if faster rates of queue growth should occur, a corresponding increase in the metering rate must be enacted to return the average queue length at or below the length specified by $q_i(k+1) = q_i(k) + \Delta T(d_{i,N} - r_{i,N})$. Thus, we have the constraints

$$\Delta q_{i,MIN}(k) \leq \Delta q_i(k) \leq \Delta q_{i,MAX}(k) \quad \forall i, k \quad \text{Eqn. 6- 22}$$

for the PC-RT problem, where $\Delta q_{i,MIN}(k)$ is defined as

$$\Delta q_{i,MIN}(k) = -q_{i,N}(k) \quad \forall i, k \quad \text{Eqn. 6- 23}$$

indicating that the queue cannot become negative, and $\Delta q_{i,MAX}(k)$ is defined as

$$\begin{aligned} \Delta q_{i,MAX}(k) &= q_{i,MAX} - q_{i,N}(k) \quad \forall i, k \\ q_{i,MAX} &= q_{i,storage} + z_i \end{aligned} \quad \text{Eqn. 6-24}$$

where z_i is the value of the "overflow" variable for ramp i from the solution of the area-wide coordination problem. Notice that $q_{i,MIN}(k)$ and $q_{i,MAX}(k)$ are functions of time according to the nominal growth rate $q_{i,N}(k)$.

Control variable modeling

Similar to the definition of the $\Delta q_i(k)$ variables, we define the control variable $\Delta r_i(k)$ as in(6-21). For realistic representation of the problem, we cannot expect to be able to make changes to the metering rate $r_i(k)$ on time scales of 5-15 seconds, as frequently as the macroscopic model equations are re-evaluated, Thus the constraints

$$\begin{aligned} \Delta r_i(k) &= \Delta r_i(k+1) = \dots \Delta r_i(k+t) & \forall i \\ \Delta r_i(k+t+1) &= \Delta r_i(k+t+2) = \dots \Delta r_i(k+2t) \\ &\vdots \\ \Delta r_i(k+(Z-1)t+1) &= \Delta r_i(k+t+2) = \dots \Delta r_i(k+Zt) \end{aligned} \quad \text{Eqn. 6- 25}$$

are added to the formulation to specify that the ramp metering rate $r_i(k)$ can only be changed once each t re-evaluation steps, where t is the number of re-evaluation steps per minute and Z is the number of minutes in the time horizon K .

Derivation of the PC-RT cost coefficients from the QP solution

The final aspect of the mathematical description of the PC-RT rate regulation optimization problem is the derivation of the weighting coefficients c_j and $c_{i,q}$ in the objective function (6-17). These weighting coefficients indicate the relative importance

of the state variables $\Delta q_i(k)$ and $\Delta p_j(k)$ and can be used to integrate the solution results of the area-wide coordination problem with the local PC-RT optimizations. To illustrate this, consider when the upper-layer area-wide coordination optimization problem is solved, a set of dual multipliers λ_ζ and slacks ε_ζ are obtained for the constraints $g_\zeta(r, q) \geq 0 \quad \forall \zeta$ in the problem. When a constraint is *tight*, (i.e. $g_\zeta(r, q) = 0$), the dual multiplier is nonzero $\lambda_\zeta \neq 0$ and the slack is zero $\varepsilon_\zeta = 0$. When a constraint has *slack*, (i.e. $g_\zeta(r, q) > 0$) then $\lambda_\zeta = 0$ and $\varepsilon_\zeta \neq 0$. The dual multiplier of a constraint indicates the *price* the decisionmaker/modeler would pay to obtain an additional unit of that resource. In other words, the dual multiplier indicates, in terms of the units of the objective cost function being optimized, how much cost would be incurred if the RHS of the constraint was increased by one unit. Therefore, when a constraint has slack, $\varepsilon_\zeta \neq 0$, the stakeholder should *not* be willing to “pay” (or incur cost) to obtain an additional unit of that resource.

Computational procedure to obtain cost coefficients

Consider that the area-wide QP is a maximization problem. If all constraints are converted to $g_\zeta(r, q) \leq 0$ constraints, then $\lambda_\zeta \geq 0 \quad \forall \zeta$. In the following derivation, we do not distinguish between c_j and $c_{i,q}$, terming each cost parameter generically c_ζ . Let λ_{min} be the smallest non-zero dual price such that

$$\lambda_{MIN} = \min_{\zeta} \{ \lambda_\zeta \mid \lambda_\zeta > 0 \} \quad \text{Eqn. 6- 24}$$

and let λ_{MAX} be the largest non-zero dual price such that

$$\lambda_{MAX} = \max_{\zeta} \{ \lambda_\zeta \mid \lambda_\zeta > 0 \}. \quad \text{Eqn. 6- 25}$$

Constraints with $\lambda_\zeta \neq 0$ are then assigned the cost coefficient

$$c_\zeta = \frac{\lambda_\zeta}{\lambda_{MAX}} + 1 \quad \text{Eqn. 6- 26}$$

to scale the cost coefficient relative to the maximum cost. Constraints with zero dual prices $\lambda_\zeta = 0$ are assigned the cost coefficient

$$c_\zeta = \frac{\lambda_{MIN}}{\lambda_{MAX}} \left(\frac{cap_\zeta - \varepsilon_\zeta}{cap_\zeta} \right) + 1 \quad \text{Eqn. 6- 27}$$

where cap_ζ is the capacity (right-hand-side) of constraint ζ and ε_ζ is the slack of constraint ζ . In this manner, the constraints with zero dual prices are linearly scaled to the cost coefficient of the constraint with λ_{MIN} . In the event that *no* constraints have non-zero dual multipliers, $\{\lambda_\zeta \mid \lambda_\zeta > 0\} = \emptyset$, then the cost coefficients are assigned according to

$$c_\zeta = \left(\frac{cap_\zeta - \varepsilon_\zeta}{1 - \varepsilon_{MIN}} \right) + 1 \quad \text{Eqn. 6- 28}$$

where ε_{MIN} is the lowest slack ratio (i.e. the closest constraint to being tight as a percentage of its capacity) such that

$$\varepsilon_{MIN} = \min_\zeta \left\{ \frac{cap_\zeta - \varepsilon_\zeta}{cap_\zeta} \right\}. \quad \text{Eqn. 6- 29}$$

This pricing method could be considered a "trickle-down" approach to reflect the congestion conditions at the interchange and the area-wide coordination priorities in each of the local PC-RT optimization problems. Thus, the combination of the dual multipliers λ_ζ and the slack values ε_ζ incorporates information from all of the following: cost coefficients of the area-wide QP, queue storage size at each ramp, demand to each ramp, and capacity of each freeway segment.

Modification to the linearization procedure for unstable conditions

When a section of the freeway is in the *unstable* regime of the characteristic equation $V_j = \rho_j \cdot v_j$ such that $\rho_j > \rho_{j,crit}$, the linearized system (6-13) does not effectively describe the system dynamics. Recall from Chapter 1, Figure 1-1, the shape of the volume-density

characteristic. Note that in Figure 1-1, the axes are reversed from the development to follow. When the density ρ_j in a section is overcritical, a backward-traveling (i.e. upstream-traveling) congestion wave is created with wave speed equal to the slope of the volume-density curve at the congested point [Newell, 1993]. Linearizing about the stable equilibrium point $\rho_{j,N}$ results in approximating the *upstream* traveling wave with a *downstream* traveling wave as shown in Figure 6-2. Here, the grey points indicate the current state, nominal state, and approximated state. Notice that the uppermost point is the estimate of the current state by linearizing the system about the nominal state and the predicted volume is significantly higher than the volume-density characteristic function, and unrealistic.

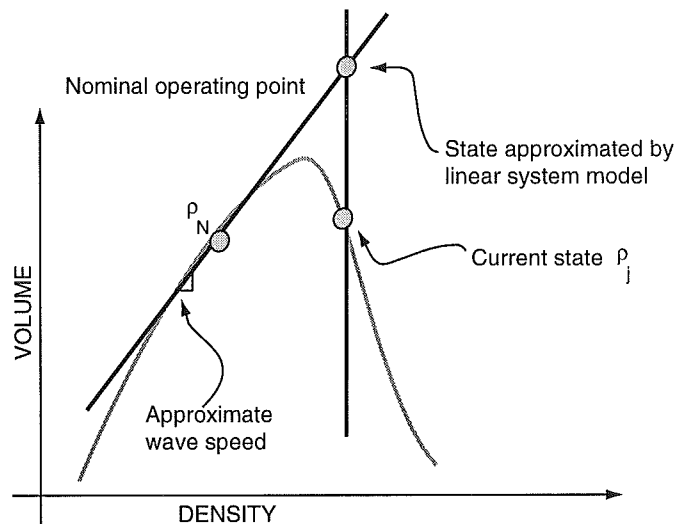


Figure 6- 2. Example of incorrect wave-speed model for congested section

One approach to remedy this modeling inaccuracy is to re-define the nominal point ρ_N in the congested section at the current *congested* density $\tilde{\rho}_N = \rho_j$ and solve the PC-RT control problem linearized about ρ_j . As shown in previous research, this model results in containment of the congestion to that segment when appropriately strong feedback control is applied [Payne et al., 1985]. However, the congestion cannot be eliminated if the cost function (i.e. LQR quadratic feedback rules) penalizes *both* positive and negative deviations from the nominal point. With such a cost function, the feedback rule increases

upstream metering rates $r_{i-t}(k) \geq r_{i-t,N} \mid t > 0$ when the density decreases in the congested downstream segment $\rho_j(k) \leq \tilde{\rho}_{j,N}$ (and thus $V_j(k) > \tilde{V}_{j,N}$). Therefore the downstream segment is kept in the congested regime $\rho_j(k) \rightarrow \tilde{\rho}_{j,N} \mid \tilde{\rho}_{j,N} \geq \rho_{crit}$. Obviously this is not an acceptable solution to the congestion problem.

Congestion *can* be eliminated in the oversaturated segment by re-solving the upper-layer area-wide coordination problem with a reduced maximum flow rate $\tilde{V}_{j,N} < V_{j,MAX}$ in the congested section. Simply put, the congestion in that segment cannot be eliminated unless $V_{j-1}(k) < V_{j+1}(k)$ for the time period until $\rho_j(k) < \rho_{j,crit}$ in the congested section [Papageorgiou, 1983]. This results in nominal points for upstream segments that must be below the maximum volume level of the oversaturated segment, that is $\tilde{V}_{j-t,N} \leq \tilde{V}_{j,N} \mid t > 0$, since in our "steady-state" assumption all vehicles that enter the system will exit the system during the time horizon T . Otherwise, if $\tilde{V}_{j-t}(k) \geq \tilde{V}_{j,N} \mid t > 0$, vehicles will be stored in section j , and possibly sections $j-t \mid t > 0$.

An example of this case is illustrated in Figure 6-3. Here, a small five-section freeway is illustrated with segments A-E. Segment E is congested, (i.e. $\rho_E(k) > \rho_{E,crit}$) at point E in Figure 6-3. Correspondingly, in the solution to the area-wide coordination problem, the flow rates at the uncongested (i.e. $\rho_j(k) \leq \rho_{j,crit}$) segments A, B, C, and D must be less than the flow rate at E, indicated by points A, B, C, and D in Figure 6-3.

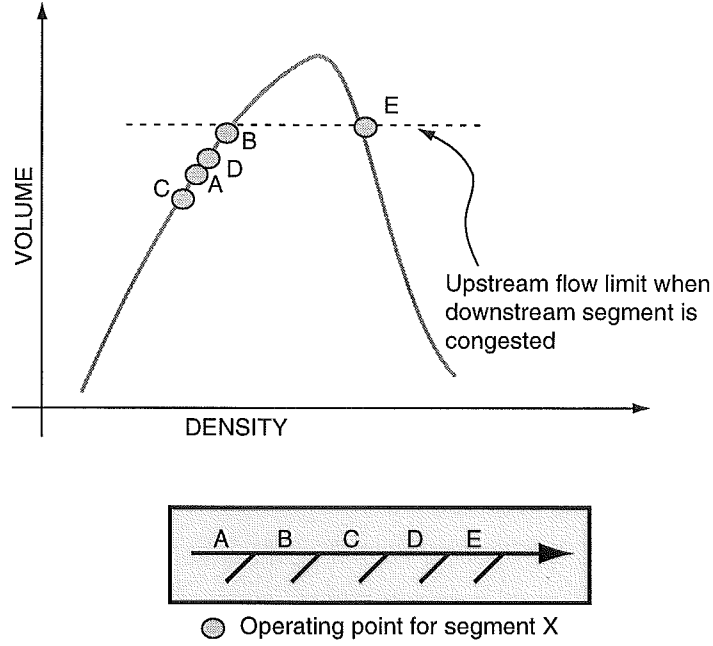


Figure 6- 3. Overcapacity segment results in upstream area-wide flow limitations

The condition illustrated in Figure 6-4 can result in congestion clearing by turning “off” any traffic-responsive control and using the nominal metering rates $r_{i,N}$ from the area-wide *capacity-limited* coordination problem. However, it seems apparent that some additional performance benefit may be realized by using a real-time problem that is appropriately structured in the congested section. Namely, an optimization problem formulation that does not result in forcing the state to the undesirable congested point.

Hence, to consider the application of the PC-RT rate regulation optimization problem to the congested section, recall that (6-16) is a *linear* cost function. Thus, the cost function already reflects the fact that $\rho_j(k+1) > \rho_{j,N}$ should be penalized when $\rho_j(k) > \rho_{j,crit}$ but $\rho_j(k+1) < \rho_{j,N}$ should not. Second, the *nominal* density $\rho_{j,N}$ in the over-capacity section is re-defined as the *critical* density $\rho_{j,crit}$ at which the maximum flow $V_{j,MAX}$ occurs. The system is linearized, however, at the currently congested point $\rho_j(k)$. The resulting approximation is illustrated in Figure 6-4.

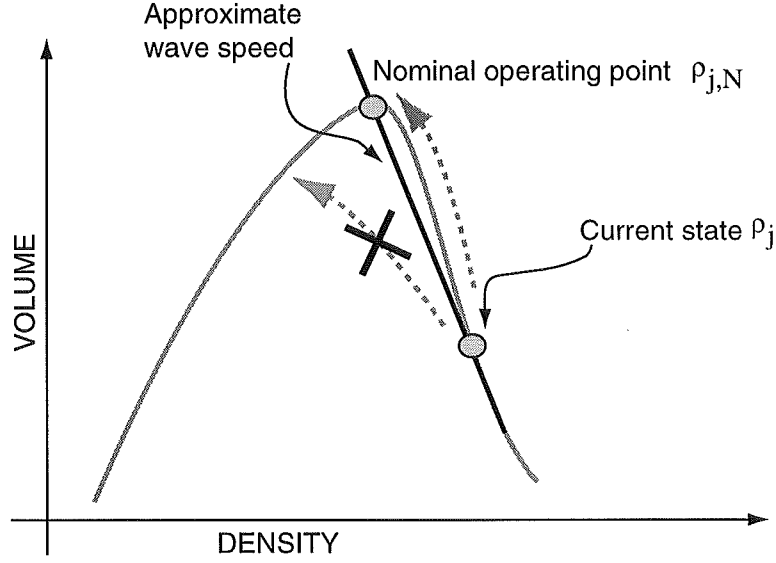


Figure 6- 4. Alternative model for the over-capacity situation

This model assumes that, in a dynamical sense, the combination of flow and density in a congested section must travel *continuously* along the fundamental curve from the congested point $\rho_j(k)$ to the maximum flow point $\rho_{j,crit} = \rho_{j,N}$, as shown in Figure 6-5, assuming infinitesimally small time intervals, Δt . The density does not jump from point to point discontinuously as uncongested flow resumes, as illustrated in Figure 6-4 by the “X-ed” arrow from the congested regime to the uncongested regime. The assumption of continuous flow change is further substantiated by Figures 1-2 and 1-3 from the Phoenix area freeway indicating that the transition from congested to uncongested flow is continuous during the rush-hour period.

Hence, for the specification of the PC-RT optimization problem in the congested section, constraints (6-14) must be modified in the case where $\rho_j > \rho_{j,crit}$ since we must periodically re-linearize the model in the congested segment at the current point as the state $\rho_j(k)$ approaches the critical value $\rho_{j,crit}$. The state constraints are thus given as

$$\max[\rho_{j,MIN}, \rho_{j,crit} - \rho_j] \leq \Delta\rho_j(k) \leq \min[\rho_{j,jam} - \rho_j, \rho_{j,MAX}] \quad \forall j \quad \text{Eqn. 6- 30}$$

replacing constraints (6-14) for the congested section only. The control constraints are unchanged from the previous definition of (6-15). This periodic re-linearization and resolution of the area-wide QP, and therefore re-formulation of the PC-RT optimization problem, is illustrated in Figure 6-5. In this figure, the topmost diagram shows the initial condition when the freeway segment is severely congested. Hence, the first approximation to the (backward-traveling) wave speed is shown as the line between the nominal point and the current point. After a few minutes, the constraints on the upstream flow rates will begin to lower the congested density and, correspondingly the volume exiting the congested segment increases. Thus, in the middle diagram, the wave-speed is re-approximated and the area-wide coordination problem is re-solved with the new (higher) volume in the congested segment. Similarly, as the density continues to decrease and the volume continues to increase in the congested segment, the PC-RT and area-wide coordination problems are re-formulated and re-solved again, as shown in the final diagram at the bottom of Figure 6-5.

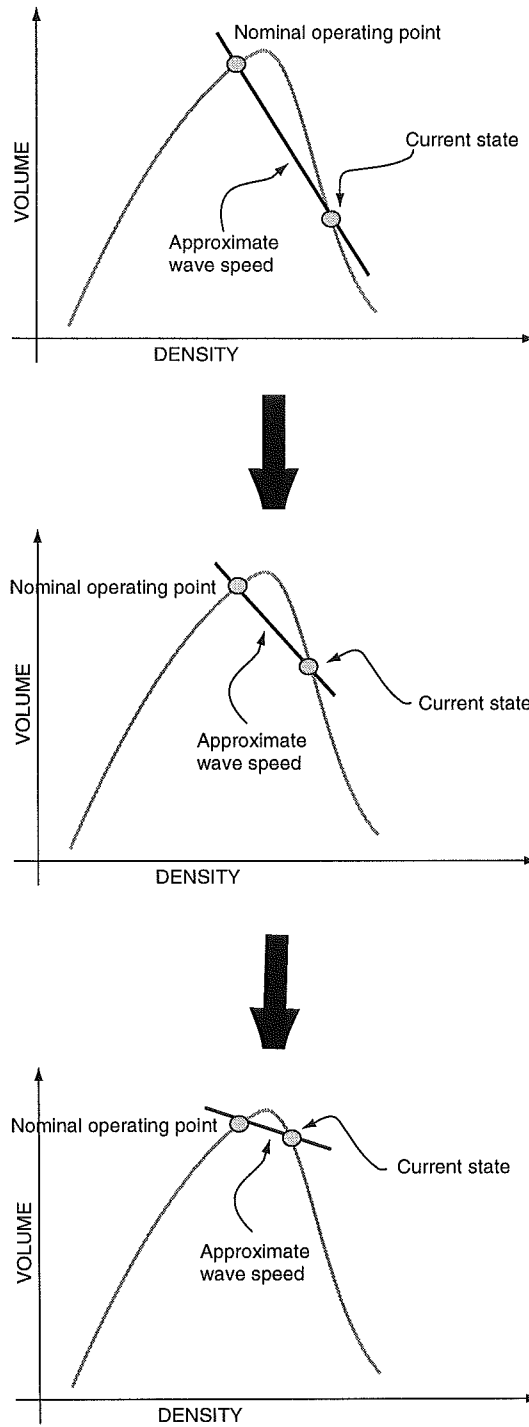


Figure 6- 5. Re-linearization for PC-RT and periodic solution of the QP

Summary of PC-RT mathematical formulation

To summarize the development of the PC-RT optimization problem to this point, we have specified the following single-objective *monolithic* control problem for the entire freeway system to save additional travel time where possible

$$\min \sum_{k=1}^K \left[\sum_{j=1}^J c_j \Gamma_j \Delta \rho_j(k) + \sum_{i=1}^M c_{i,q} \Delta q_i(k) \right] \quad \text{Eqn. 6- 31}$$

st

$$\Delta \rho(k + \Delta T) = \Delta \rho(k) + \Delta T(\mathbf{A}^* \Delta \rho + \mathbf{B}^* \Delta r) \quad \forall k$$

$$\Delta q_i(k + 1) = \Delta q_i(k) + \Delta T[\Delta d_i(k) - \Delta r_i(k)] \quad \forall i, k$$

$$\text{if } \rho_j(k) \leq \rho_{j,crit} \quad \forall j$$

$$\max[\rho_{j,MIN}, -\rho_{j,N}] \leq \Delta \rho_j(k) \leq \min[\rho_{j,crit} - \rho_{j,N}, \rho_{j,MAX}] \quad \forall k$$

otherwise

$$\max[\rho_{j,MIN}, \rho_{j,crit} - \rho_j] \leq \Delta \rho_j(k) \leq \min[\rho_{j,jam} - \rho_j, \rho_{j,MAX}] \quad \forall k$$

$$r_{i,MIN} - r_{i,N} \leq \Delta r_i(k) \leq \hat{r}_{i,MAX} - r_{i,N} \quad \forall i, k$$

$$\Delta q_{i,MIN}(k) \leq \Delta q_i(k) \leq \Delta q_{i,MAX}(k) \quad \forall i, k$$

$$\Delta q_{i,MIN}(k) = -q_{i,N}(k) \quad \forall i, k$$

$$\Delta q_{i,MAX}(k) = q_{i,MAX} - q_{i,N}(k) \quad \forall i, k$$

$$q_{i,MAX} = Q_i + z_i \quad \forall i$$

$$\Delta r_i(k) = \Delta r_i(k + 1) = \dots \Delta r_i(k + t) \quad \forall i$$

$$\Delta r_i(k + t + 1) = \Delta r_i(k + t + 2) = \dots \Delta r_i(k + 2t)$$

\vdots

$$\Delta r_i(k + (Z - 1)t + 1) = \Delta r_i(k + t + 2) = \dots \Delta r_i(k + Zt)$$

This formulation is a linear programming problem with unrestricted-sign variables $\Delta r_i(k)$, $\Delta q_i(k)$, and $\Delta \rho_j(k)$.

This LP has $J*K$ freeway variables $\Delta \rho_j(k)$ and $M*K$ queue variables $\Delta q_i(k)$, $3*M*K$ queue constraints and $4*J*K$ freeway constraints. The PC-RT rate regulation optimization problem is thus appreciably larger than the upper-layer area-wide QP as a *single* optimization problem, but not insurmountably large, in the "large-scale" sense, given

modern computing power and solution algorithms. The difficulty, however, of solving this monolithic optimization problem (and the primary reason for the hierarchical decomposition of the freeway management problem) is that the solution depends on the predictions of (1) $\Delta d_i(k)$ and (2) $\Delta V_{o,IN}(k)$ to the freeway system at the ramp meter and upstream freeway input, respectively, over the short-term time horizon $K\Delta T$. The following development establishes the reasoning behind the decomposition of the monolithic control problem into subproblems.

Difficulty in solving the monolithic PC-RT optimization problem

Consider first that, over the time horizon $K\Delta T$ of the PC-RT optimization problem, the nominal inputs, as assumed in the area-wide coordination problem, $d_{i,N}(k) = \bar{d}_{i,N}$ and $V_{o,IN}(k) = \bar{V}_{o,IN}$ remain constant. At $k=0$ we measure the current d_i and $V_{o,IN}$ and find that $V_{o,IN}(k) \neq \bar{V}_{o,IN}$ and/or $d_i(k) \neq \bar{d}_{i,N}$. At this point we try to predict what *could* happen in the next 5-7 minutes to $d_i(k)$ and $V_{o,IN}(k)$ and the corresponding effect on the total system.

Consider each ramp having three “fundamental” predicted scenarios, for $d_i(k)$:

- (1) increases $d_i(k) > d_i(k-1)$,
- (2) decreases $d_i(k) < d_i(k-1)$, or
- (3) remains constant $d_i(k) = d_i(k-1)$.

Similarly the upstream freeway input $V_{o,IN}(k)$ can

- (1) increase further $V_{o,IN}(k) > V_{o,IN}(k-1)$,
- (2) decrease $V_{o,IN}(k) < V_{o,IN}(k-1)$, or
- (3) remain constant $V_{o,IN}(k) = V_{o,IN}(k-1)$.

Thus, in eqn. 6-44 we need to evaluate 3^{M+1} LPs to evaluate the optimal rate combinations for each *possible* combination of predicted inputs at the M ramps and 1 freeway upstream input. After solving 3^{M+1} LPs and measuring the actual demands at each ramp in the next observation period, we must search (and therefore also store) a $M+1$ -dimensional table to extract the appropriate set of metering rates.

For a typical problem size, say 10 ramps, evaluating just three predicted scenarios at each ramp necessitates solving $3^{11} = 177,147$ LPs each minute to obtain the new metering rate for each possible combination of predicted demands and be guaranteed that the “optimal” metering rate is applied in the next minute. Thus, our approach is to use a sub-optimal procedure for the solution of the PC-RT rate regulation problem and *decompose* the monolithic problem into M smaller problems. For example, in an area-wide coordination problem with 10 ramps, we evaluate just $(3*3)*M = 90$ LPs and search M 2-dimensional rate tables for the appropriate rate for the next observation period. Such a scheme is more applicable for real-time control, hopefully with little degradation in performance.

An argument to suggest that the performance would *not* degrade considerably is simply to consider the fact that the predictions being made are *approximate*, and thus the inaccuracy of the predictions likely outweighs any benefits of real-time solutions between the ramps by solving many, many inaccurate monolithic PC-RT rate regulation optimization problems. Further, because of the short-time horizon (5-7 minutes) and continual re-optimization in a rolling-horizon fashion (1-2 minutes), the inaccuracies of predictions should be suitably mitigated. The next section details how the monolithic PC-RT rate regulation optimization problem is decomposed into sub-problems.

Decomposition of full optimization problem into subproblems

For continuous linear systems, many methods exist to decompose a large-scale system into subsystems to satisfy various objectives of reliability, bounded stability, or strength of interaction between subsystems [Isaksen and Payne, 1973; Sandell et al., 1978; Goldstein and Kumar, 1982; Papageorgiou and Mayr, 1982]. For the PC-RT rate regulation optimization problem, our decomposition relies on the concept of the “strength of interaction” between subsystems as well as the requirement to keep the optimization problem *computationally tractable* in real-time. That is, all computations must be completed in less than one-minute for all ramps in the system.

Hence, although we are guided by the analytical approaches of previous work, we implement a straightforward decomposition principle of including *one* ramp per

subsystem, with enough upstream freeway to satisfy the predicted travel-time requirement and enough downstream freeway to allow for anticipatory “reactivity” to downstream congestion. In general, from the topology of most freeway networks (i.e. most ramps are spaced approximately 1 mile apart), our decomposition scheme results in the following *overlapping* subsystem description [Kahng et al., 1984; Haimes et al., 1990] shown in Figure 6-6.

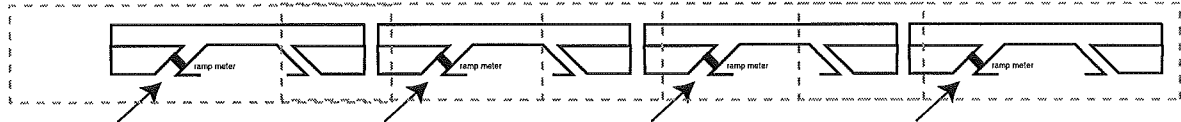


Figure 6- 6. Typical overlapping subsystem decomposition

Hence, we *decouple* (6-33) into M individual-ramp subproblems of the same form

$$\min \sum_{k=1}^K \left[\sum_{m=j-t}^{j+1} c_m \Gamma_m \Delta \rho_m(k) + c_{i,q} \Delta q_i(k) \right] \quad \text{Eqn. 6- 32}$$

st

$$\Delta \rho_m(k + \Delta T) = \Delta \rho_m(k) + \Delta T (\mathbf{A}^{*j} \Delta \rho_m(k) + \mathbf{B}^{*j} IO_j(k)) \quad \forall m \in I^j, \forall k$$

$$IO_j(k) = [V_{j-B,IN}(k), 0, 0, \dots, 0, r_i(k), V_{j+1,out}(k)]$$

$$\Delta q_i(k + 1) = \Delta q_i(k) + \Delta T [\Delta d_i(k) - \Delta r_i(k)] \quad \forall k$$

$$\text{if } \rho_m(k) \leq \rho_{m,crit} \quad \forall m \in I^j$$

$$\max[\rho_{m,MIN}, -\rho_{m,N}] \leq \Delta \rho_m(k) \leq \min[\rho_{m,crit} - \rho_{m,N}, \rho_{m,MAX}] \quad \forall k$$

otherwise

$$\max[\rho_{m,MIN}, \rho_{m,crit} - \rho_m] \leq \Delta \rho_m(k) \leq \min[\rho_{m,jam} - \rho_m, \rho_{m,MAX}] \quad \forall k$$

$$r_{MIN} - r_{i,N} \leq \Delta r_i(k) \leq \hat{r}_{MAX} - r_{i,N} \quad \forall k$$

$$\Delta q_{i,MIN}(k) \leq \Delta q_i(k) \leq \Delta q_{i,MAX}(k) \quad \forall k$$

$$\Delta q_{i,MIN}(k) = -q_{i,N}(k) \quad \forall k$$

$$\Delta q_{i,MAX}(k) = q_{i,MAX} - q_{i,N}(k) \quad \forall k$$

$$q_{i,MAX} = Q_i + z_i$$

$$\begin{aligned}
\Delta r_i(k) &= \Delta r_i(k+1) = \dots \Delta r_i(k+t) & \forall i \\
\Delta r_i(k+t+1) &= \Delta r_i(k+t+2) = \dots \Delta r_i(k+2t) \\
&\vdots \\
\Delta r_i(k+(Z-1)t+1) &= \Delta r_i(k+t+2) = \dots \Delta r_i(k+Zt)
\end{aligned}$$

such that A^{*j} is defined as

$$\mathbf{A}^* = \begin{bmatrix} a_{j-B,j-B} & a_{j-B+1,j-B} & 0 & \dots & 0 \\ a_{j-B,j-B+1} & a_{j-B+1,j-B+1} & a_{j-B+2,j-B+1} & & \vdots \\ 0 & & \dots & & 0 \\ \vdots & & & a_{j-1,j} & a_{j,j} & a_{j+1,j} \\ 0 & \dots & 0 & a_{j,j+1} & a_{j+1,j+1} \end{bmatrix} \quad \text{Eqn. 6- 33}$$

and B^{*j} is

$$\mathbf{B}^* = \begin{bmatrix} 1 & & & & \\ & 0 & & & \\ & & \ddots & & \\ & & & 0 & \\ & & & & 1 \\ & & & & & 1 \end{bmatrix} \quad \text{Eqn. 6- 34}$$

using a numbering convention where the most upstream density section in the freeway is designated as ρ_1 . $IO_j(k)$ is a vector of the inputs and outputs to this ramp subsystem, such as the upstream freeway input $V_{j-B,IN}(k)$, ramp input $r_i(k)$ and off-ramp output $s_i(k)$, and downstream freeway output flow $V_{j+1,OUT}(k)$. The example shown in (6-34) indicates no off-ramp in the section with the on-ramp, and no other on-ramps in the sub-system.

According to the modeling convention of the macroscopic freeway model of Chapter 4, each freeway *section* j is composed of num_j *segments* and thus the total number of segment in the freeway model is $J * num_j$. The next modeling convention is that each on-ramp or off-ramp is found in the *first* segment of a new section, such that the changes in the system are found at the 1st, num_1+1 , $num_1 + num_2 + 1, \dots, \sum_{j=1}^{N-1} num_{j-1} + 1$ segments of the whole system. Thus, the set I^j is defined as

$$I^j = [j - B, j - (B - 1), \dots, j, j + 1] \quad \text{Eqn. 6- 35}$$

where j is the segment containing the i^{th} ramp ($num_1 + num_2 + \dots + num_{(j-1)} + 1$) and B is the minimum between the number of segments upstream that constitute more than one-minute travel time, but less than two minutes, or the number of segments until another ramp is encountered. The choice of one-minute travel time is from the specification that the metering rate only be changed in one-minute intervals. Note also that as the travel-time in each segment of the system changes, as during incident conditions, the size definitions of each sub-problem may change.

This criterion for B also ensures that only *one* segment with an on-ramp occurs in each sub-problem and thus the metering rate decisions made by each PC-RT subproblem are independent. It may be necessary, however, in some real-world freeway systems where there are several ramps less than one-minute of travel-time apart, to consider their real-time optimization together using a different decomposition scheme. This possibility is not explored in this research.

We also assume, in the specification of the subproblems, that the volume immediately downstream $V_{j+1,OUT}$ of the subsystem remains at the nominal flow rate $\bar{V}_{j+1,OUT,N}$ during the PC-RT optimization horizon. This is a simplifying assumption since the downstream flow could *also* be considered a location for detection of flow anomalies $V_{j+1,OUT}(k) \gg V_{j+1,OUT,N}$ or $V_{j+1,OUT}(k) \ll V_{j+1,OUT,N}$. With this addition, however, the number of optimization problems would increase three-fold, and hence we would need to solve for $3 \times 3 \times 3 = 27$ prediction scenarios. This possibility may be studied in future research.

Scenario prediction

To provide the *pro-active* capability of the PC-RT optimization layer we need to predict the possible changes to the ramp demands and the upstream freeway flow for each sub-problem. When a significant deviation to the flow rate is detected one minute upstream of the ramp section $V_{j-B}(k) \notin [V_{j-B,N} \pm \Delta V_{j-B,N}]$ on the freeway and/or at the ramp meter

entrance $d_i(k) \notin [d_{i,N} \pm \Delta d_{i,N}]$ (also one-minute upstream of the interchange flows), a *set* of possible future scenarios is created, each with a corresponding LP optimization problem. Since the detected flow variation is one minute upstream, we do not need to select the appropriate rate immediately. Thus, the procedure is as follows

- (a) Detect the anomaly $d_i(k) \notin [d_{i,N} \pm \Delta d_{i,N}]$, and/or $V_{j-B}(k) \notin [V_{j-B,N} \pm \Delta V_{j-B,N}]$,
- (b) create the set of predictions and solve for travel-time minimizing metering rates with the PC-RT subproblem decomposition and store in a table,
- (c) in the next minute, sample the flow rate $V_{j-B}(k+1)$ and $d_i(k+1)$ again one-minute upstream and evaluate which predicted scenario was closest to being realized,
- (d) apply the appropriate metering rate to “react” to the realized system state from the stored table of possible rates for ramp i , and
- (e) repeat the entire process for the current measurement.

Hence, the predictions do not necessarily need to be entirely accurate, but they do need to reflect both

- (a) the *range* of conditions possible over the short-term time horizon of 5-7 minutes, and
- (b) the *possible dynamics* over the short-term time horizon of 5-7 minutes.

In addition, at this level of decoupling, we neglect the fact that, for sub-problems in the *interior* of the system, we have more upstream information than just one-minute of travel-time. Neglecting this information does not seem egregious since the uncertainty of the entire flow process and the fact that our predictions are approximate (at best) will probably outweigh any benefits of using the additional upstream flow information for flow predictions in the interior of the freeway.

As previously stated, to represent the range of possible future states three *fundamental* prediction scenarios are created for the two input flows to each sub-problem: *increasing* flow linearly, *decreasing* flow linearly, and *remaining* at the current level, as illustrated is

Figure 6-7. The three predicted scenarios for the upstream freeway segment are labeled A_{fy} , B_{fy} , and C_{fy} and A_{rd} , B_{rd} , C_{rd} identify the three predicted demand scenarios to the ramp meter.

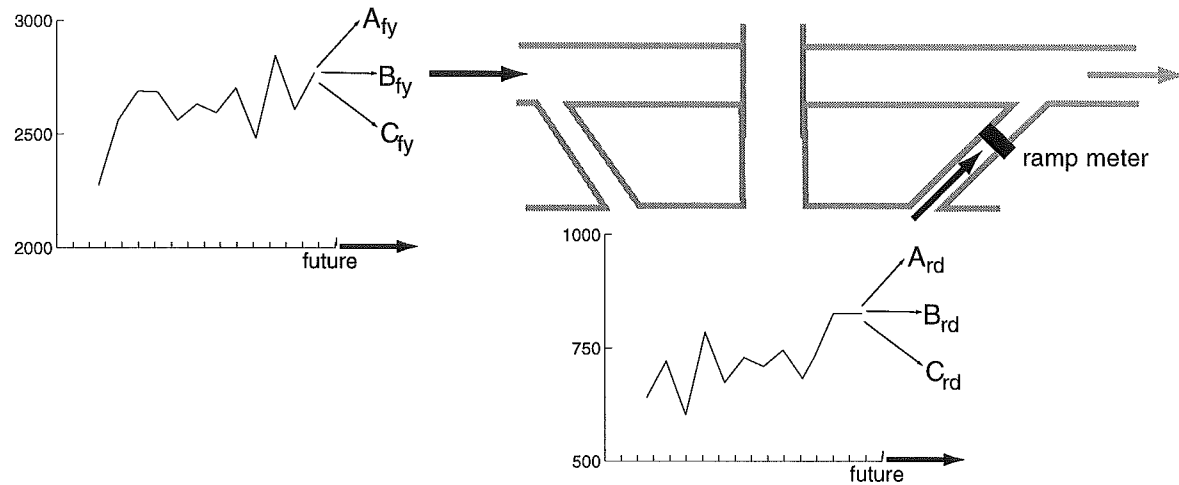


Figure 6- 7. Predicted trends for a given subproblem

Scenario prediction example

Recall that we plan the PC-RT optimization of a sub-problem *only* when an anomaly $d_i(k) \notin [d_{i,N} \pm \Delta d_{i,N}]$ and/or $V_{j-B}(k) \notin [V_{j-B,N} \pm \Delta V_{j-B,N}]$ is detected at one or both of the input streams, so the "remaining at current level" prediction could or could not be an anomalous condition. For example, consider a subsystem where the nominal demand to the ramp meter is 600 veh/hr and a demand of 750 veh/hr is detected one minute upstream. At the same subsystem the nominal demand to the upstream freeway segment is 1700 veh/hr-lane and a flow of 1850 veh/hr-lane is detected. Assume that the volume level at the upstream freeway input $V_{j-B,IN}(k) = 1850$ constitutes an anomalous flow condition for this subsystem. When this is detected, three predictions are created for the possible ramp flows and three predictions are created for the upstream freeway volume over the next five to seven minutes, for a total of *nine* possible scenarios. We do not

require that both the freeway and ramp demands be anomalous to begin PC-RT optimization.

Consider also, for example, that ramp demands can change, on average, (after smoothing out short-term detector errors and removing the cyclical effect of the interchange traffic signal) approximately 15 veh/hr per minute, either linearly positive or negative, such that

$$d_i(k+1) = d_i(k) \pm 15. \quad \text{Eqn. 6- 36}$$

In addition, assume that on average (after smoothing) freeway volume can change approximately 60 veh/hr per minute linearly, such that

$$V_{j-B,IN}(k+1) = V_{j-B,IN}(k) \pm 60. \quad \text{Eqn. 6- 37}$$

Hence, the *fundamental* predicted time-series of flows to the ramp meter for the next five minutes are given in Table 6-1 and the *fundamental* predicted time-series of flows at the upstream freeway segment for the next five minutes are given in Table 6-2.

Prediction	Minute 1	Minute 2	Minute 3	Minute 4	Minute 5
Increasing	765	780	795	810	825
Constant	750	750	750	750	750
Decreasing	735	720	705	690	675

Table 6- 1. Ramp meter demand predictions

Prediction	Minute 1	Minute 2	Minute 3	Minute 4	Minute 5
Increasing	1910	1970	2030	2090	2150
Constant	1850	1850	1850	1850	1850
Decreasing	1810	1750	1690	1630	1570

Table 6- 2. Upstream freeway demand predictions

With these predictions and the initial conditions of the subsystem (i.e. $q_i(0)$, $\rho_{j,N}$, $v_{j,N}$, $r_{i,N}$, $d_i(0)$, and $V_{j-B,IN}(0)$), nine optimization problems of the form 6-45 are solved for the metering rates over the next five minutes that minimize additional travel time in the subsystem.

Construction of scenario rate tables

The result of solving the nine optimization problems is a 3X3 table of metering rate changes from the current metering rate, as illustrated in Table 6-3. Table 6-3 indicates the metering rate modification to be applied depending on the scenario that is actually detected in the next minute at both of the input streams to the subsystem. Note also that, at each minute, the table entries will change as the initial conditions and flow rates change. Thus, for further clarification, the row and column headings *constant* in Table 6-3 do not indicate that the flow remains at the *nominal* ramp demand and/or upstream freeway flow, but rather could mean that some very high ramp demand and/or upstream freeway flow will *remain* very high for the next five to seven minutes. Thus, each time the PC-RT problem is re-solved for a subsystem, a new rate table is constructed that *only* applies to the current combination of local conditions $q_i(0)$, $\rho_{j,N}$, $v_{j,N}$, $d_i(0)$, $V_{j-B,IN}(0)$ and nominal metering rate $r_{i,N}$.

freeway demand meter demand	increasing	constant	decreasing
increasing	0	+50	+100
constant	-20	0	+20
decreasing	-50	-25	0

Table 6- 3. Rate table for next minute

Infeasible PC-RT scenarios

It may be the case, however, that some predicted flow trends result in infeasible PC-RT optimization problems. Say, for example, that the freeway density is very near the critical density $\rho_j(k) \approx \rho_{j,crit}$, and the prediction of this scenario is for increasing

upstream freeway volume $V_{j-B,IN}(k+1) = V_{j-B,IN}(k) + 60$. It may not be possible with limited control of the local ramp meter alone to keep the freeway from becoming congested $\rho_j(k) \leq \rho_{j,crit} \quad \forall j, k$ even by reducing the metering rate all the way to its minimum allowable rate $\Delta r_i = r_{i,MIN} - r_{i,N} \quad \forall k$. Recall that this minimum allowable rate is *not* the same as the *absolute* minimum allowable rate. At the same time, it may also not be possible to keep the ramp queue from spilling back (or spilling back further, if spillback was already planned by the area-wide QP) into the adjacent interchange. If either or both of these events are true, the PC-RT scenario is infeasible and the optimization routine should provide a recommendation to the upper-layer processor(s) to re-solve for new nominal metering rates if this combination of freeway and ramp flows is realized in the next minute. The existence of infeasible scenarios is demonstrated in Table 6-4.

This strategy of re-solving the area-wide coordination problem when a PC-RT problem is infeasible addresses the local nature of the PC-RT optimization problems and the hierarchical structure of the MILOS control system. First, the local PC-RT controllers do *not* know the relative congestion levels of the interchanges in the corridor and thus if the PC-RT control algorithm is allowed to make the decision to spill-back the ramp, it could well be spilling-back the queue at the *most* congested interchange in the network. In this way, the PC-RT subproblems remain sensitive to the priorities specified at the area-wide coordination layer, but in a *decentralized* manner. That is to say, for example, that a subproblem with *large* weighting coefficients c_j for the freeway variables $\Delta \rho_j(k)$ and a relatively *small* weighting coefficient $c_{i,q}$ for the ramp queue changes $\Delta q_i(k)$ does not "know" that this combination of weighting coefficients is based on the relative costs of the *entire* area-wide problem. Instead, the subproblem is provided the information that, for its part of the problem, changes to the freeway flow $\Delta \rho_j(k)$ are much more expensive to system operation than changes to the ramp queue growth $\Delta q_i(k)$. Hence, re-solving the area-wide coordination problem allows us to distribute the spill-back amongst several

upstream ramps that have lower interchange congestion levels than the current ramp, better than re-solution of the PC-RT problem with less restrictive $r_{i,MIN}$ and $r_{i,MAX}$ settings.

The second reason not to allow the PC-RT algorithm too much ability to modify the metering rate (and therefore to keep a given PC-RT scenario *feasible*) is that this flexibility could likely result in oscillatory behavior of alternately spilling-back the ramp queue (with a restrictive rate) and then “dumping” the queue at a very high metering rate. This oscillatory behavior is particularly unacceptable to practicing traffic engineers since it is viewed as defeating the purpose of ramp metering altogether [TAC, 1998].

freeway demand meter demand	increasing	constant	decreasing
increasing	resolve QP	+50	+100
constant	resolve QP	0	+20
decreasing	resolve QP	-25	0

Table 6- 4. Rate table with infeasible optimization problems

Summary

The PC-RT optimization subproblem layer of the MILOS hierarchical control system is a pro-active approach to planning real-time, traffic-responsive ramp metering rates that considers surface-street conditions in its optimization. Queue management is explicitly considered in the optimization formulation to realize “additional” travel-time savings at each ramp by taking advantage of opportunities to reduce the queue when the freeway is under-utilized and hold back additional vehicles when there is queue capacity and the freeway has a short-term surge in demand. The PC-RT optimization problem structure is tightly integrated with the solution of the area-wide coordination problem since the cost coefficients are derived from the dual multipliers and slack values of the constraints in the area-wide coordination problem. PC-RT optimization runs at each single-ramp subsystem are scheduled when a significant difference between either the upstream freeway flow rate, or the ramp demand rate (or both) is detected. A monitoring function

is required to detect such statistically-significant flow fluctuations. The theoretical foundation and typical operation of this process monitoring function, based on the concepts of statistical process control (SPC), is the subject of the next chapter.